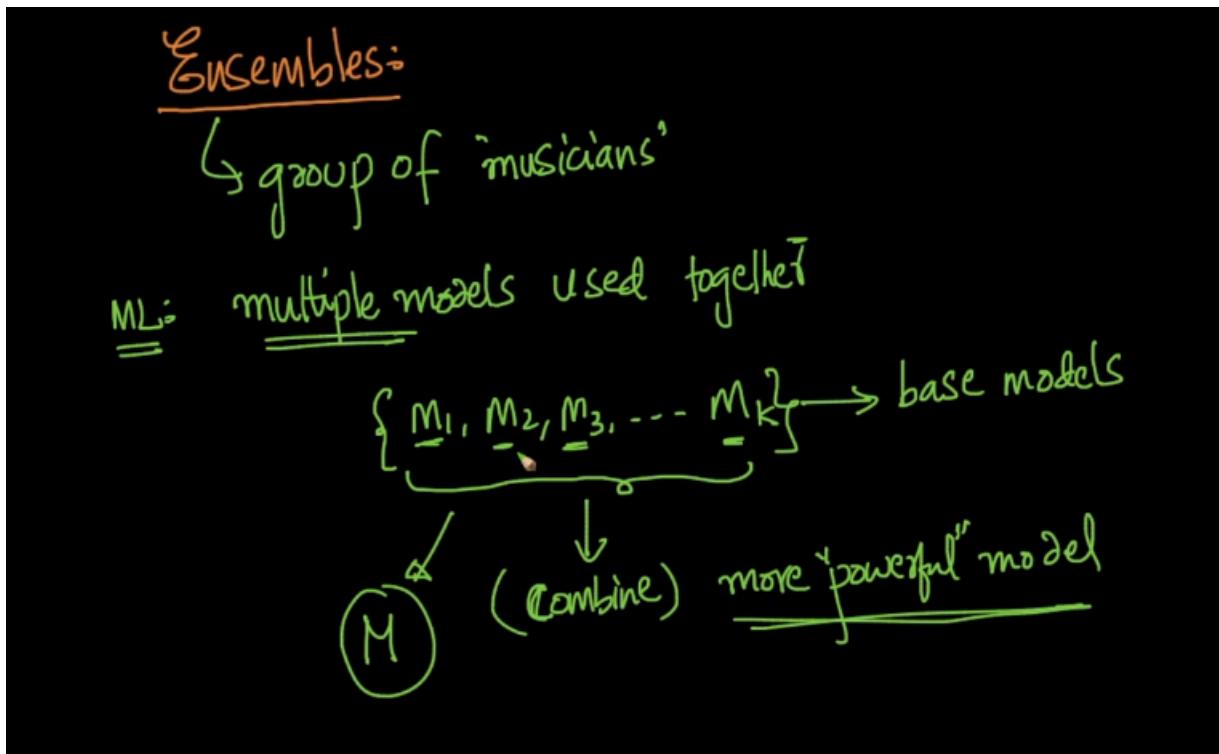


NOTES

What are ensembles?

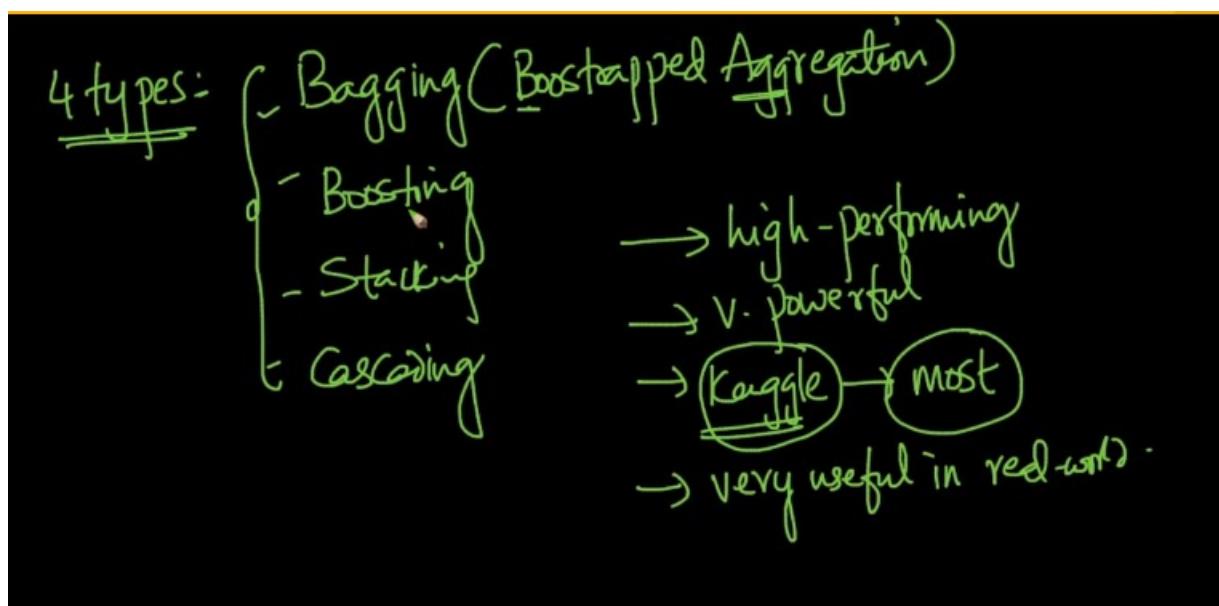
When multiple models used together to make more powerful model. Combining the models must be done ingeniously.



There are four types of ensembles:

1. Bagging. 2. Boosting. 3. stacking. 4. cascading.

The more different the models are the better we can combine them to get more powerful model.



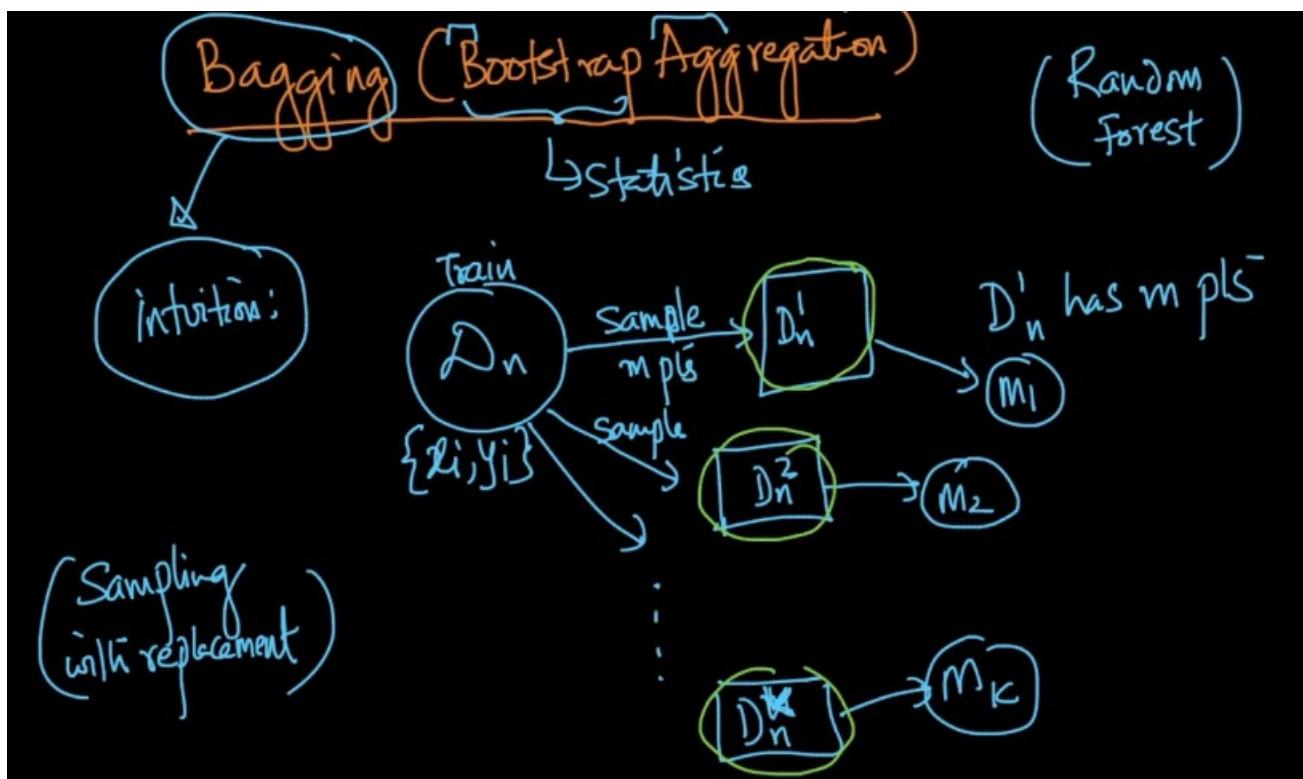
Key aspect:

$M_1, M_2, M_3, \dots, M_K$
more different These models are, }
the better you can combine them. } generic

{ Problem:- M_i :- expert

Bootstrapped Aggregation(Bagging) intuition:

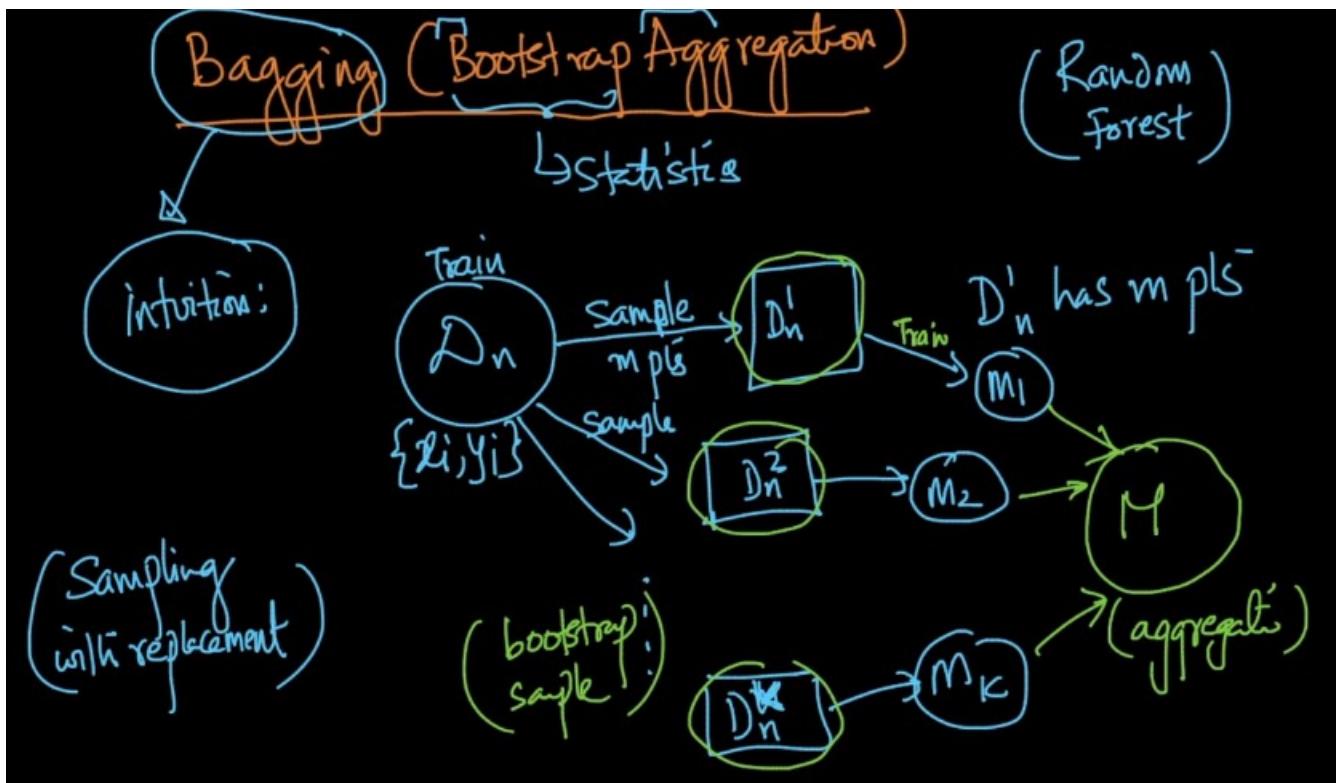
All the samples that are done sample with replacement.



Each model is built on the sample data from the complete data.

M_i is built using D_n^i of size m ($m \leq n$)
 \Rightarrow each model M_i has seen a different subset of data

After building all the models we will combine all the models.

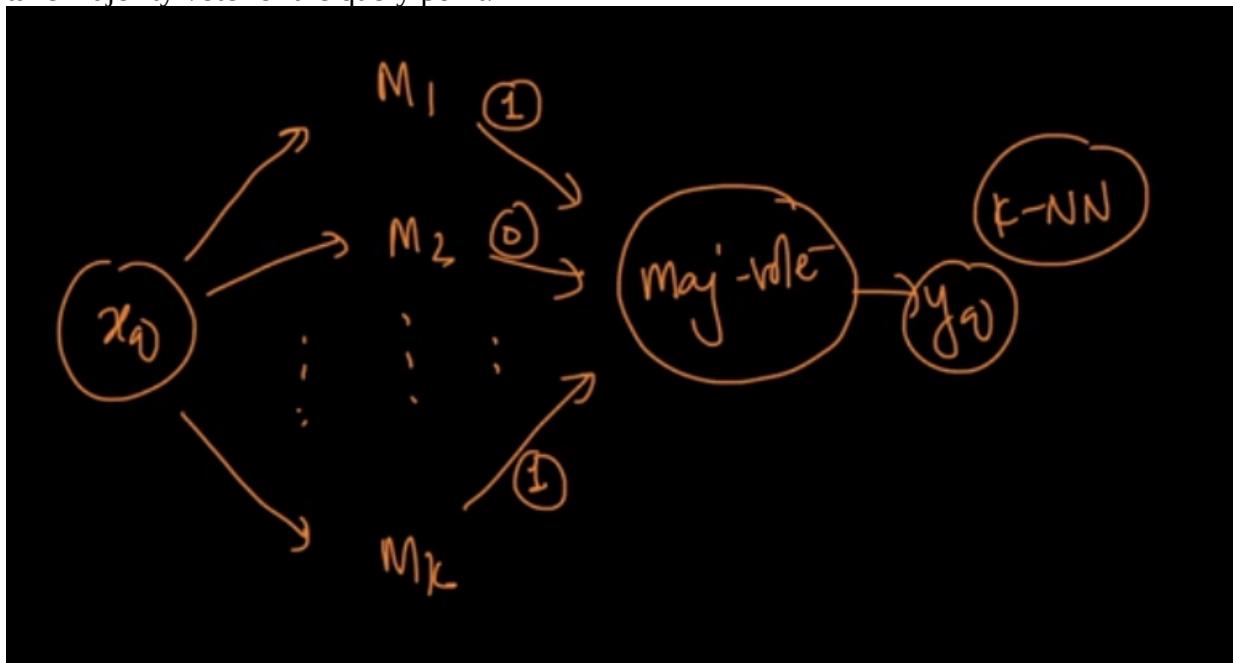


The typical way of aggregation is the majority vote in case of classification and mean/median in case of regression.

M_i is built using D_m^i of size m ($m \leq n$)
 \Rightarrow each model M_i has seen a different subset of data

Aggregation :- Classification :- Majority vote
 Regression :- mean / median

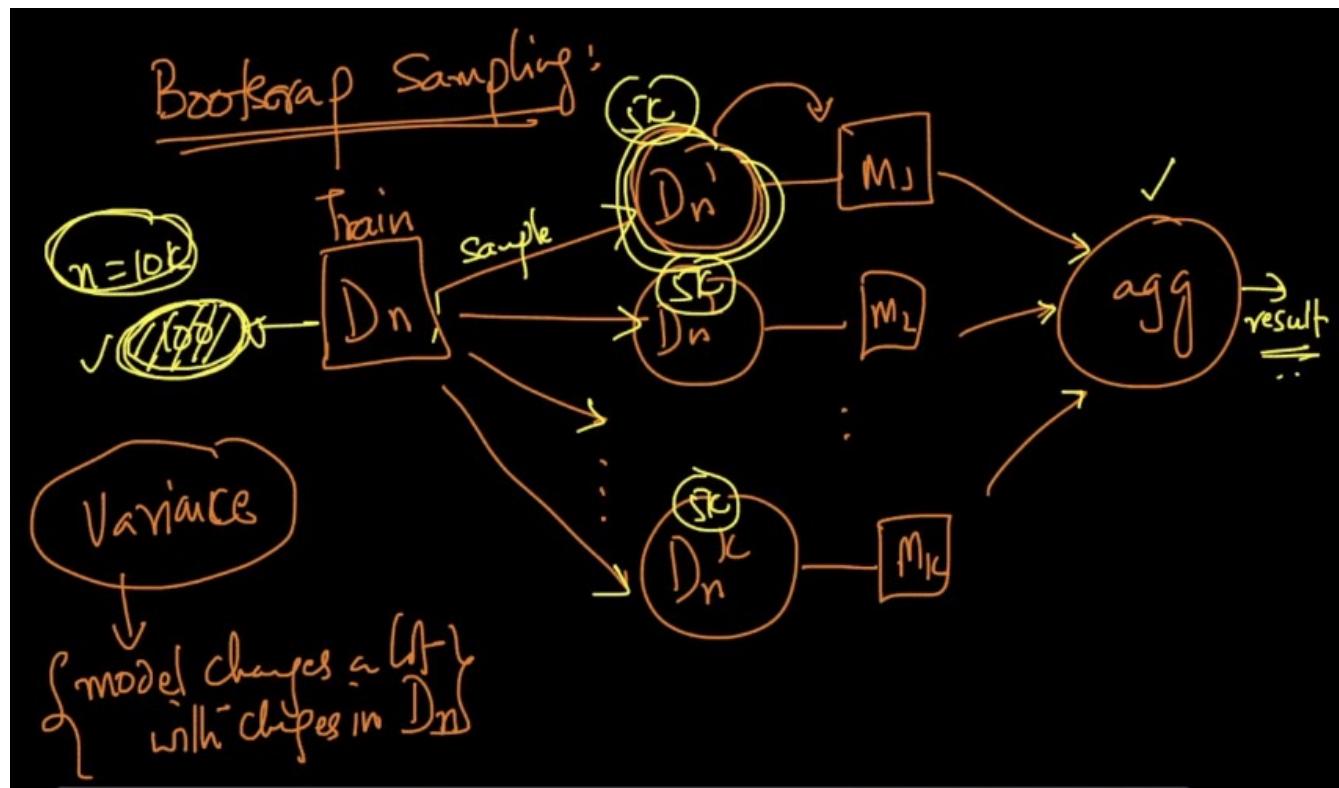
We take majority vote for the query point.



In case of regression we use the mean or median as the prediction.

By changing the original data some of the models will be changed, because every model is developed on subset of the data.

Since, we are choosing aggregation. The over all results does not change.



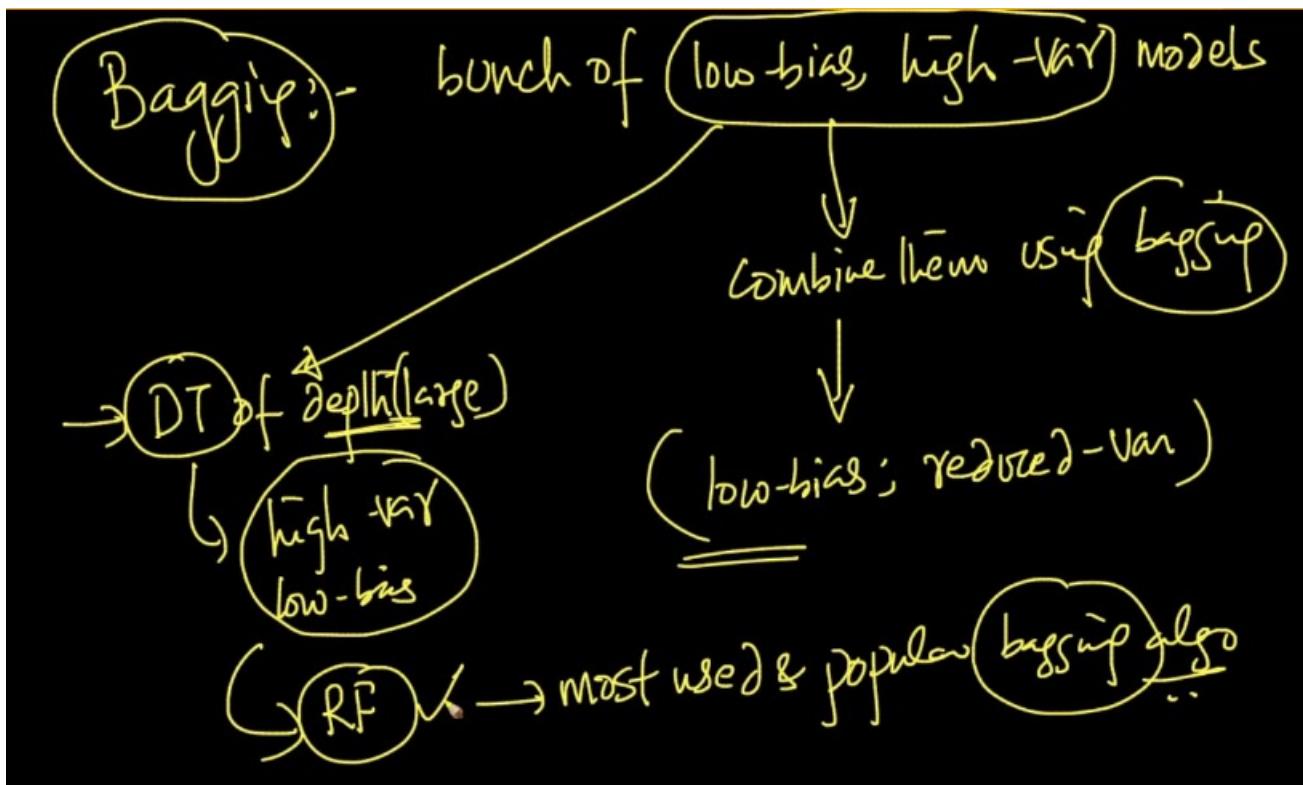
This implies bagging as the model can reduce the variance, without changing the bias.

Bagging → can reduce variance in a model without impacting the bias

$$\text{model error} = \underbrace{\text{Bias}^2}_{\text{low bias}} + \text{Var}$$

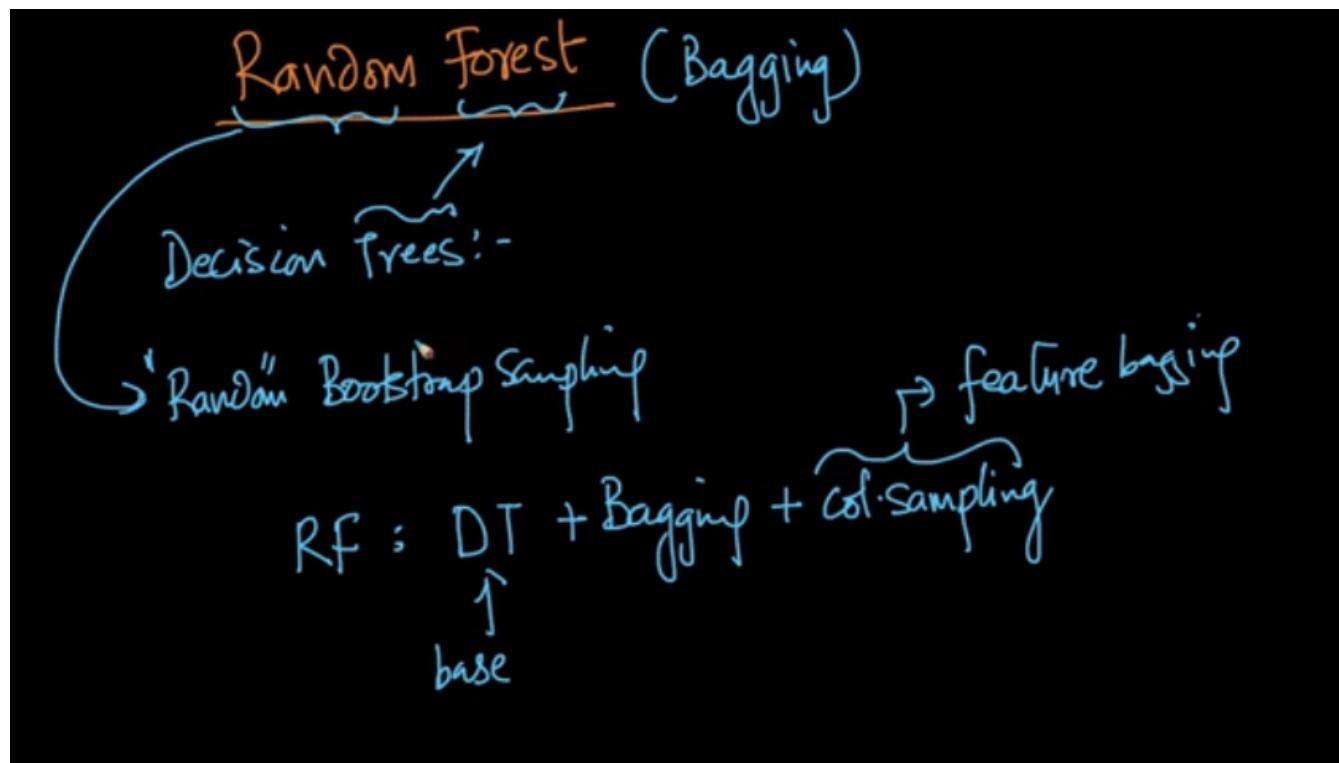
base-model (M_i): - $\underbrace{\text{low bias}}$; $\underbrace{\text{high-var model}}$
Bagging (M_i 's) → \downarrow \downarrow
low-bias; reduced variance

Characteristics of bagging:

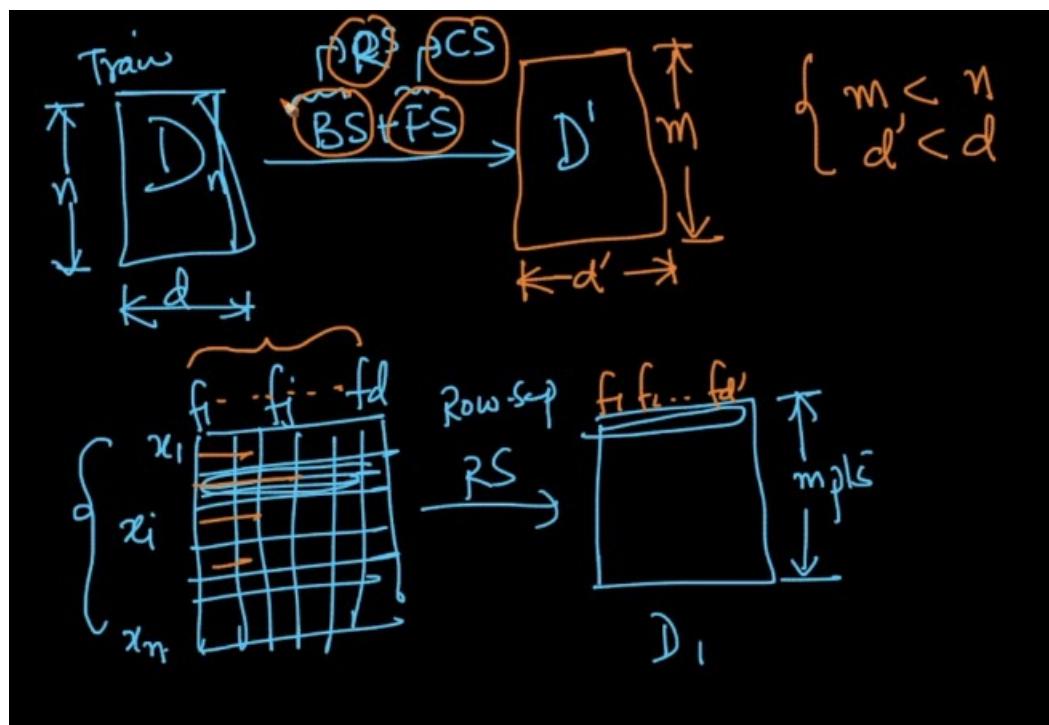


Random Forest and their construction:

It takes Decision trees as the base model for bagging and col sampling(Feature bagging).



Implementation:



The sampling of rows is done randomly with replacement.

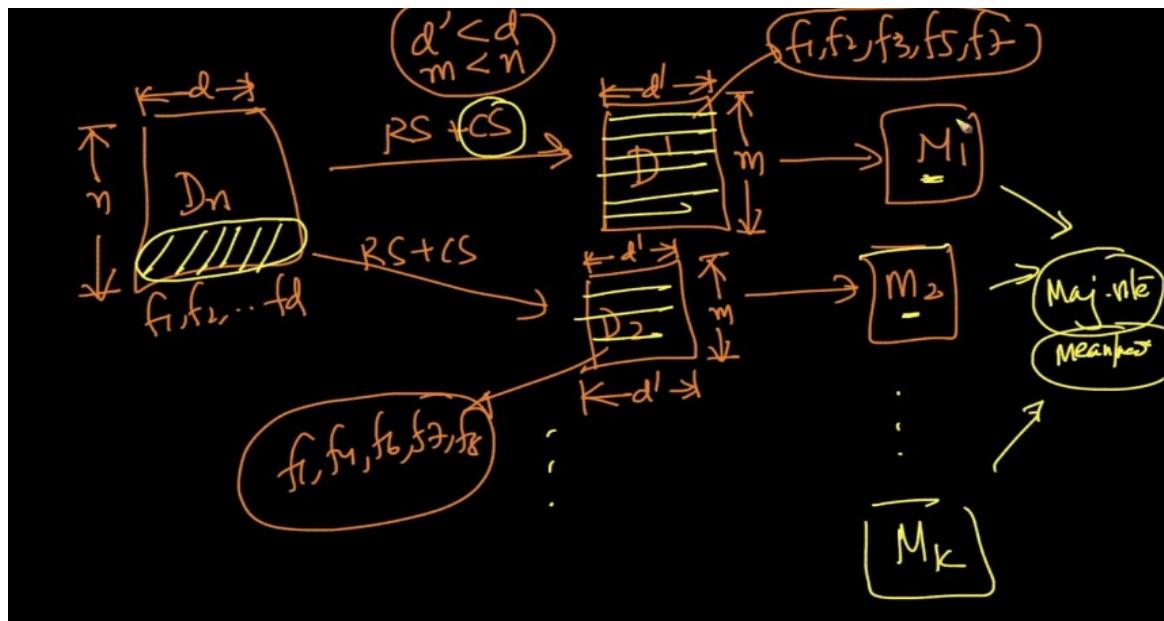
The selected sampled data set dimensions(or features) are also selected randomly.

The model M_1 is made on this data.

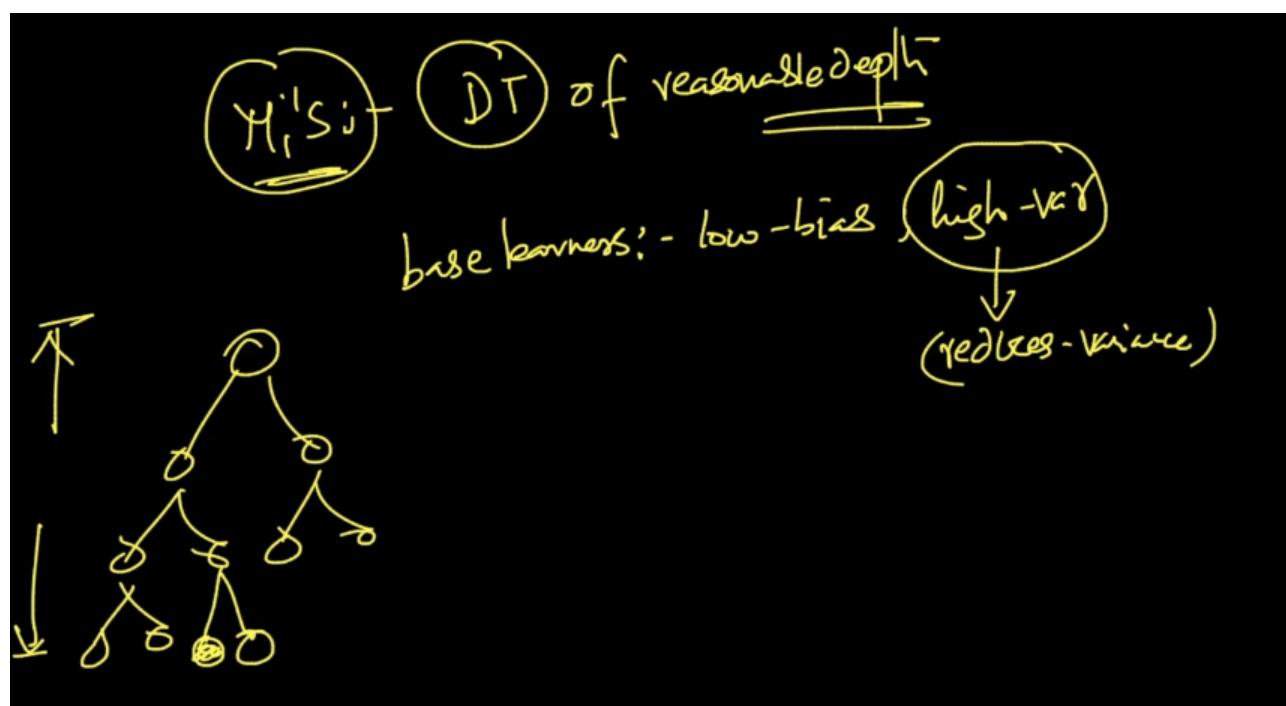
The rows and columns are also sampled randomly in Random Forest.

The above process is done for making several other models. The data sets will be more different than the bootstrap sample as the features selected are different.

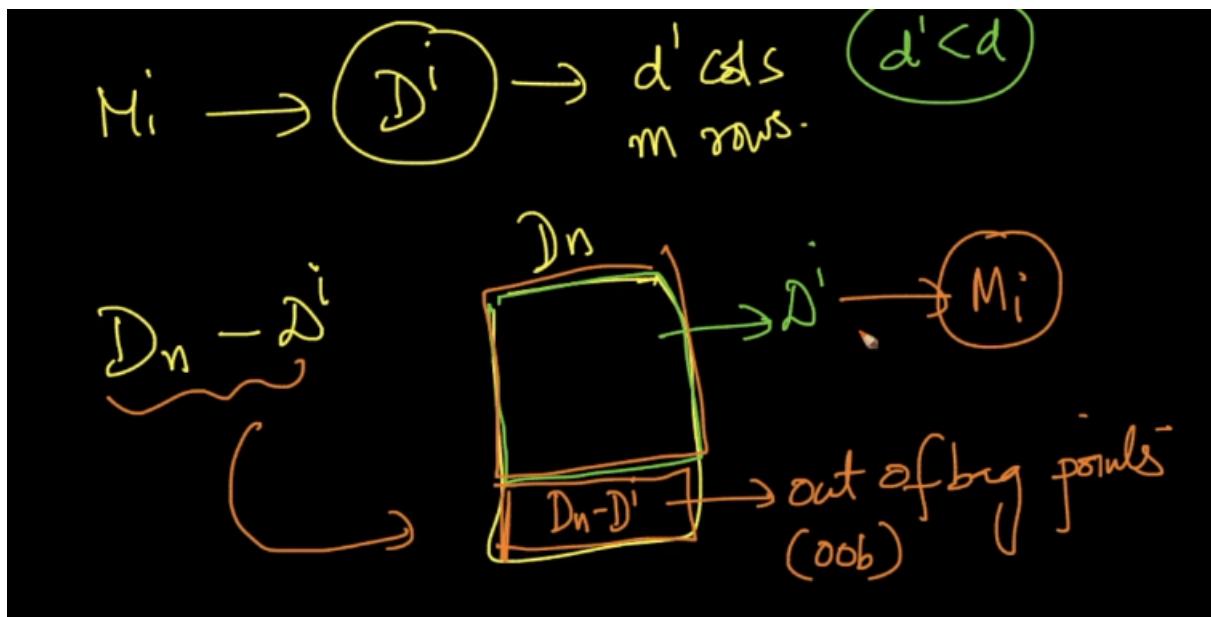
At the end we do majority vote.



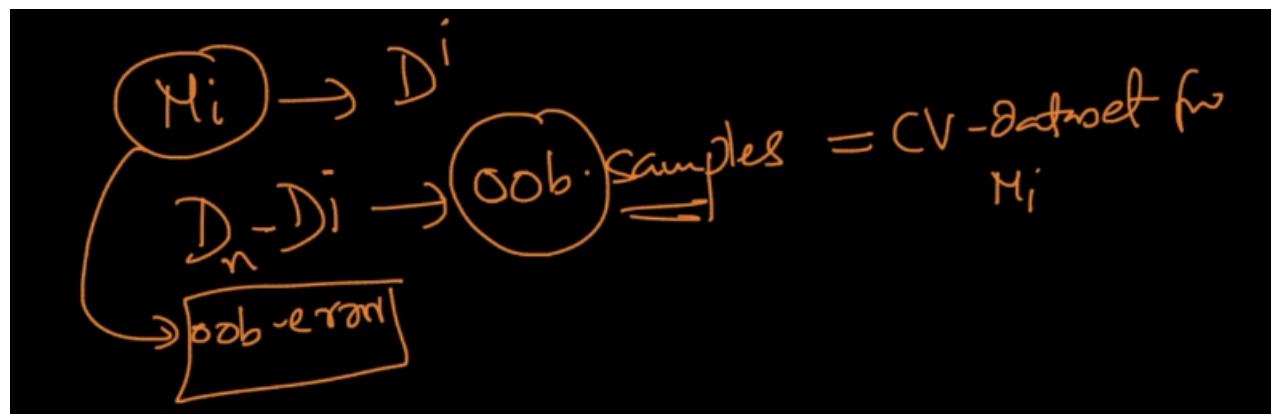
The M_i 's have no limit on the depth, making each model to over fit (high variance) does not impact the overall performance as the majority vote is done to reduce the high variance (over-fitting) of each model.



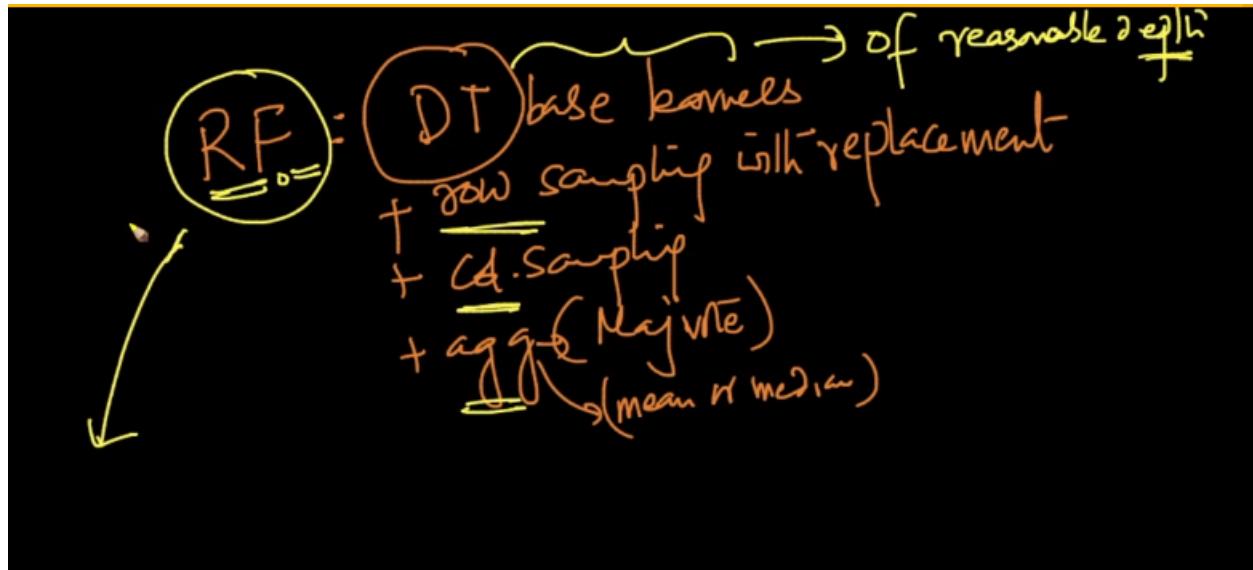
The out of bag points are nothing but the points that are not included in the sample data for training the model. These **OOB** can be used for cross-validation.



This error is called **OOB** error.

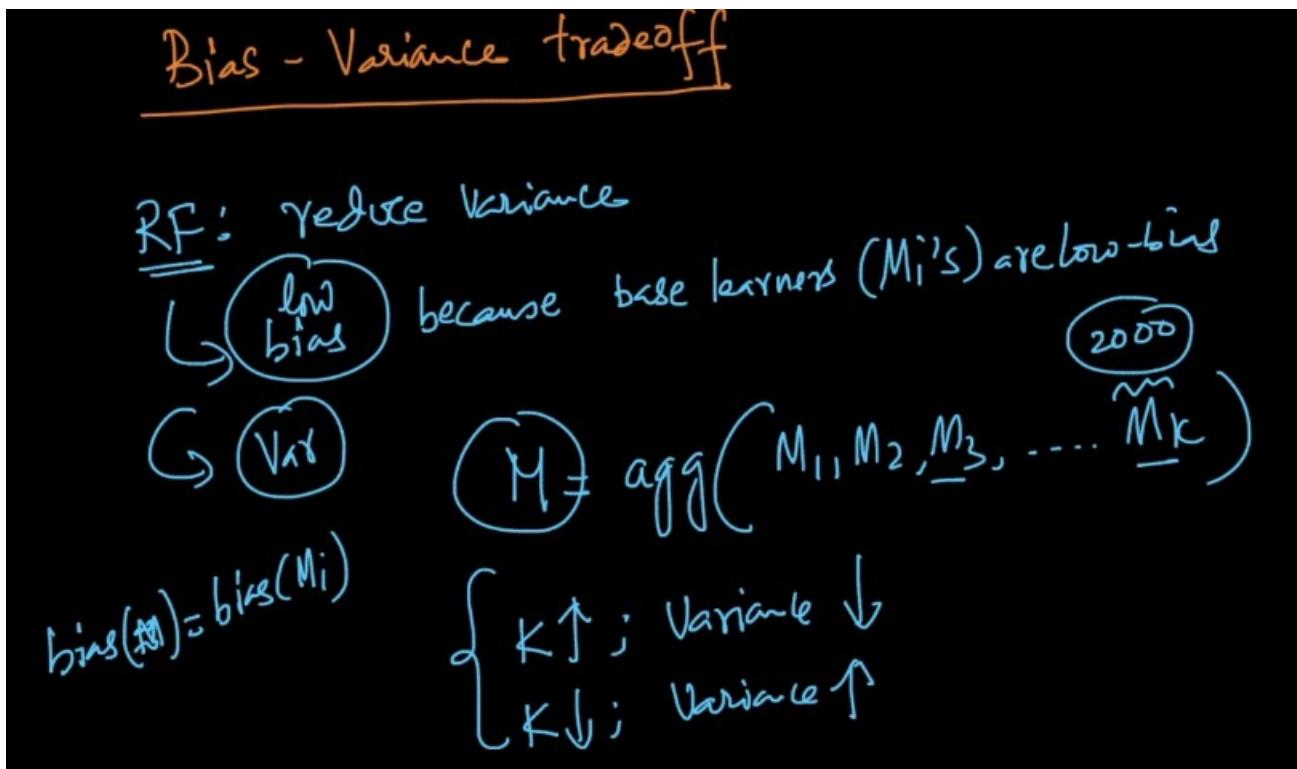


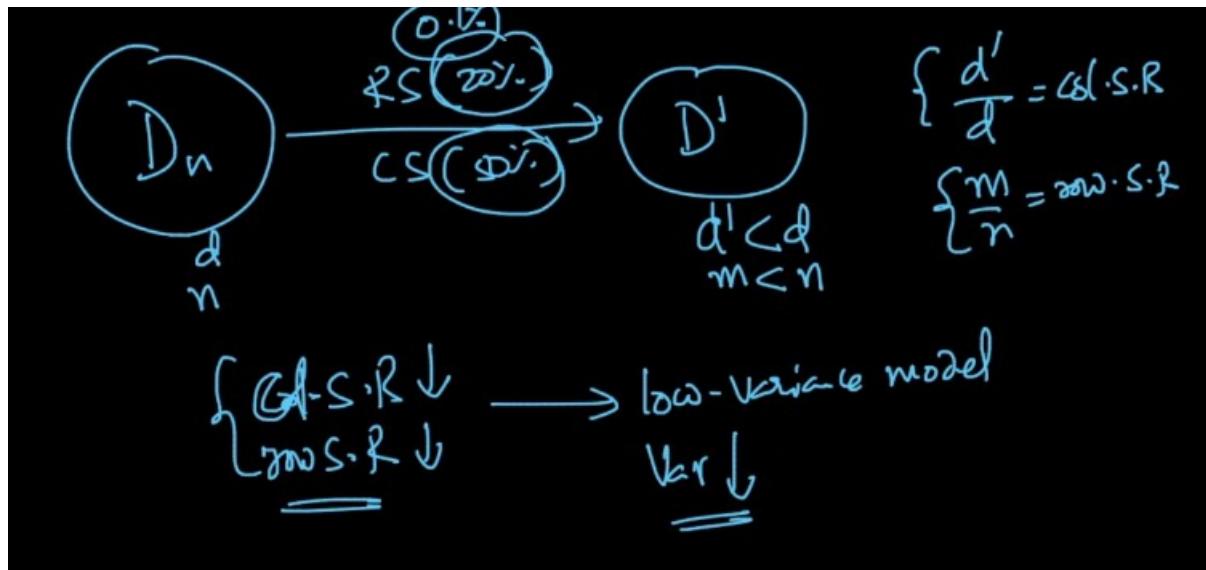
Summary:



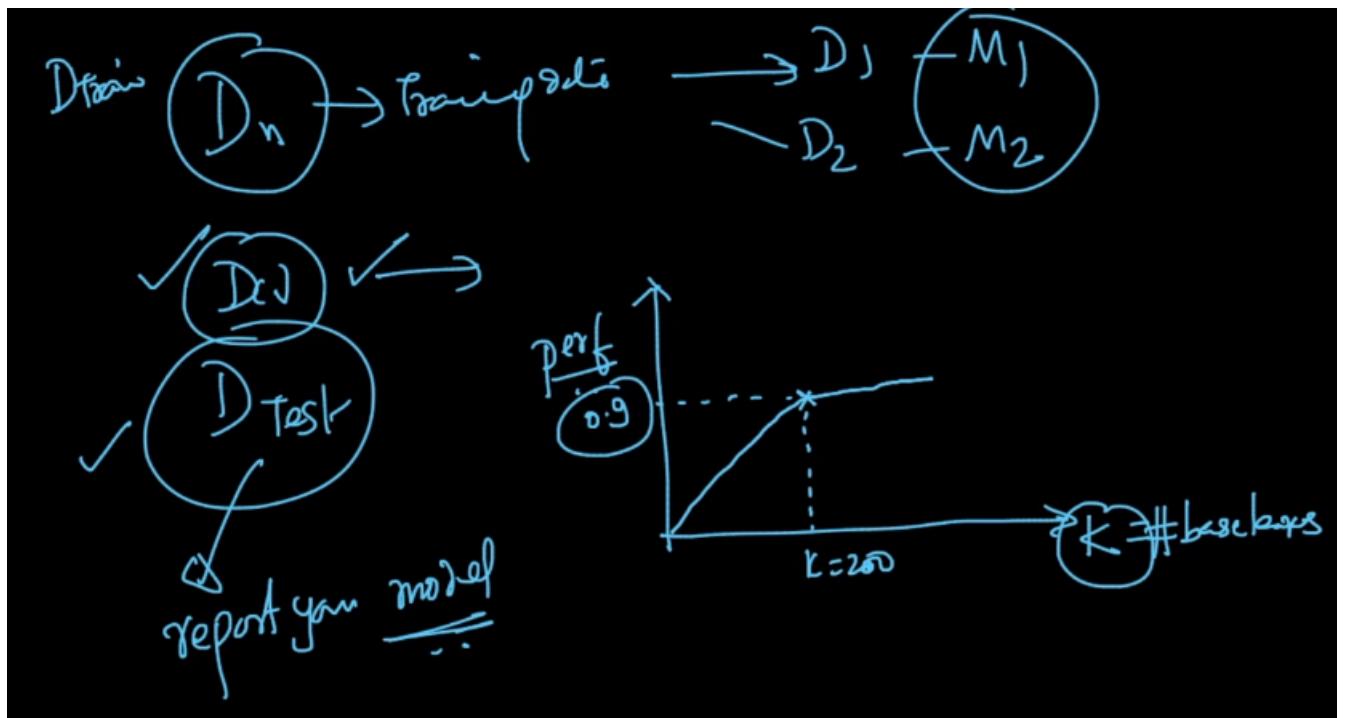
Bias – Variance trade off:

Random forest are said to have less bias because the base models are low – bias.





Often times people tend to keep the column and row sampling rate constant and keep increasing the models.

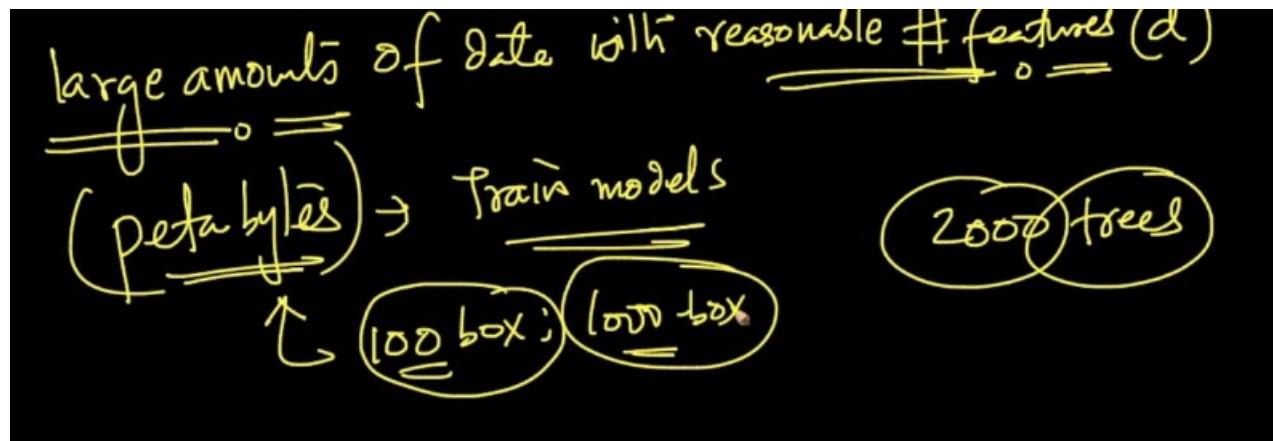
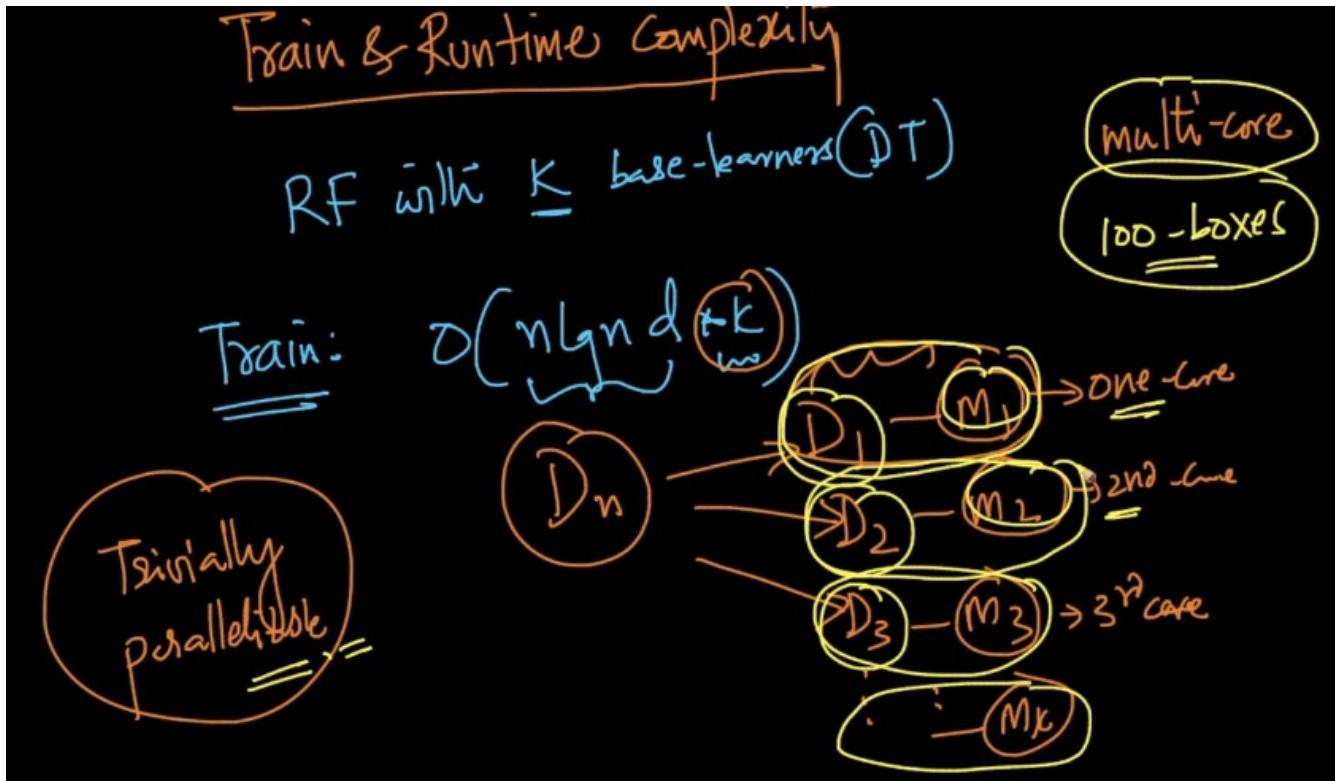


Cross-validation is used to choose the hyper parameter K (number of base learners).
The number of base learners is the important hyper parameter than the row and column rate.

Bagging: Train and Run – time Complexity:

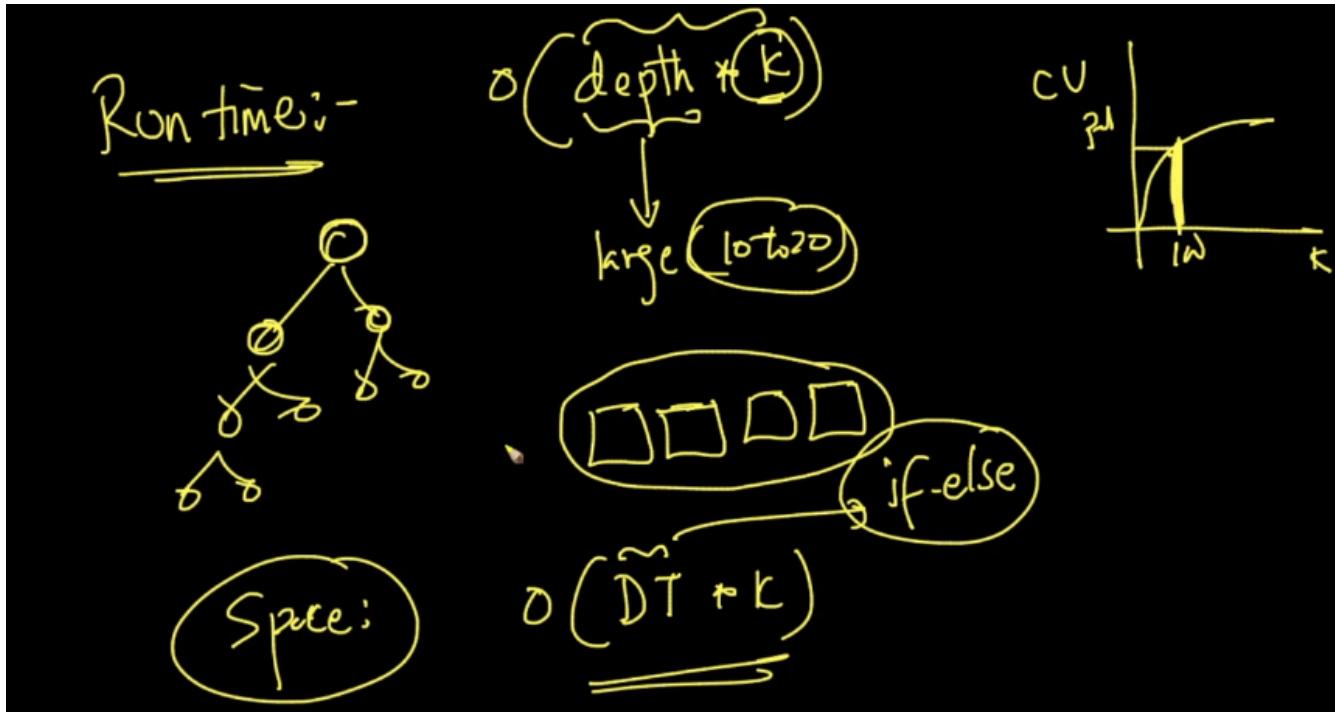
The complexity of a decision tree is: $n \log(d)$

Each of the models can be trained independently. Random forest is trivially parallelizable.



Run time complexity:

The model it self can be made into if – else condition script.

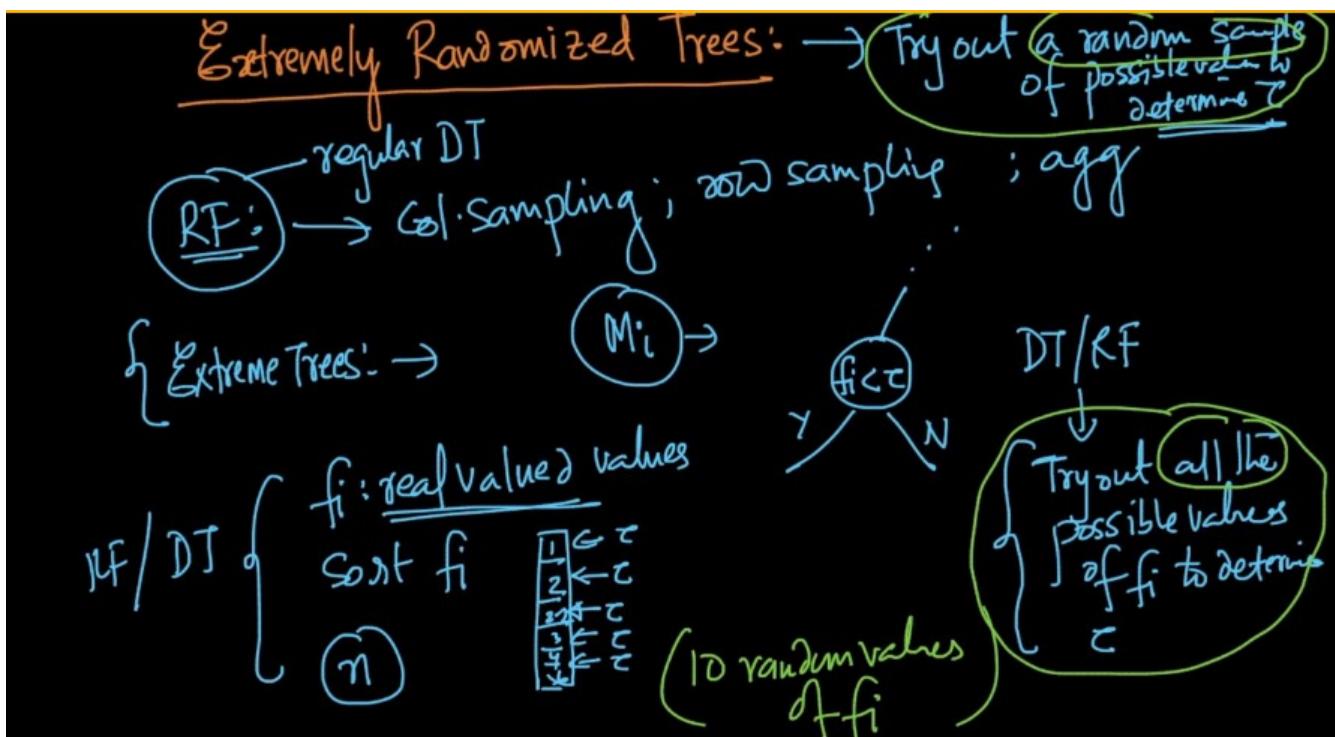


The run time and space complexity of Random forest is very low.

Extremely randomized trees:

In regular decision trees, we will try out all the values of the threshold and get the best possible value of the threshold for real – valued features, that increase the information gain.

In randomized decision trees the some values of the bag are selected randomly and then we come up with a threshold. This is our final threshold of the base model.

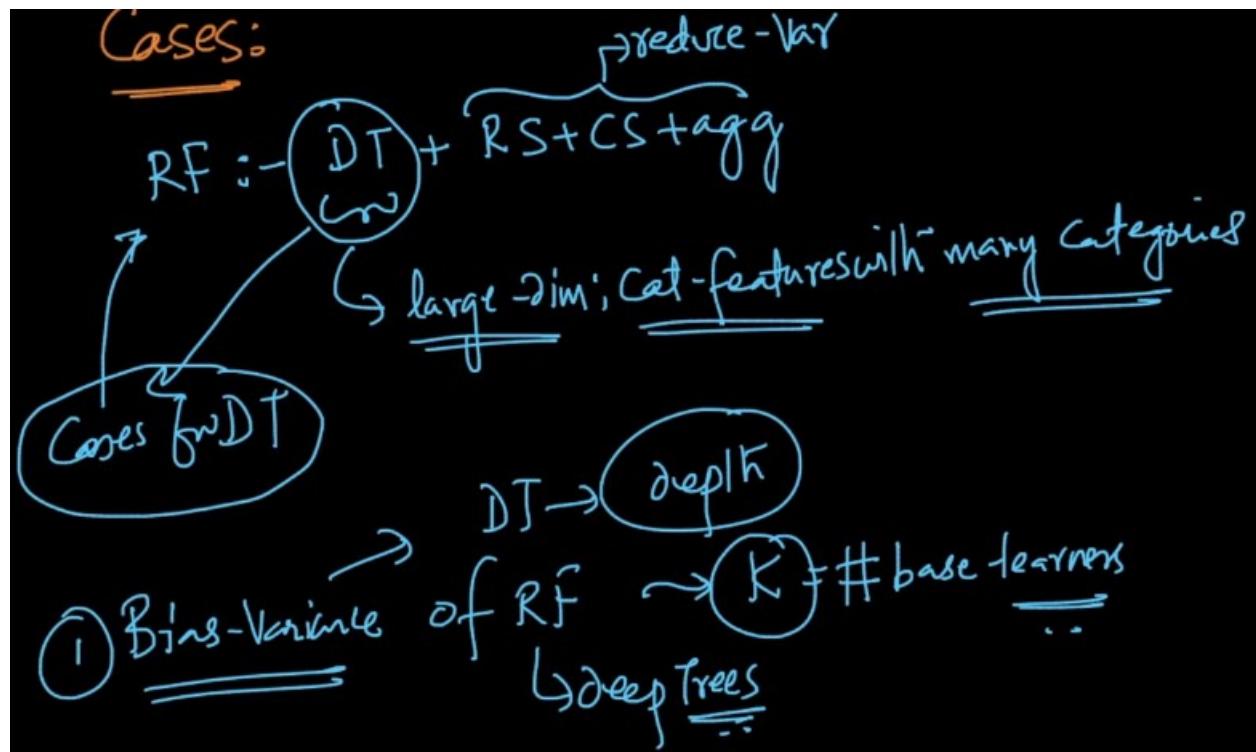


Extreme Trees: - Col Sampling + R.S + agg → RF
+ randomization when selecting c

Randomization as a way to reduce Variance

RF → C.S & R.S
 $\sum T$ → C.S + R.S + random c
 (reduce Variance better than RF)

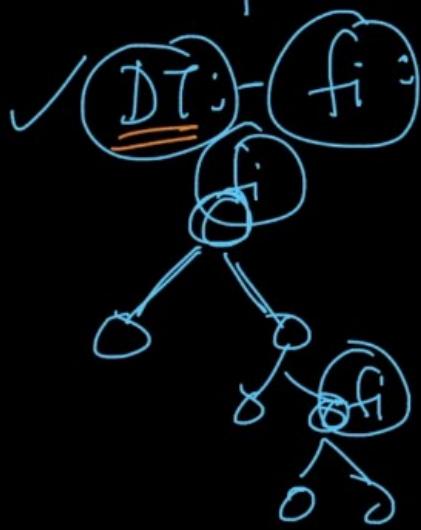
Cases:



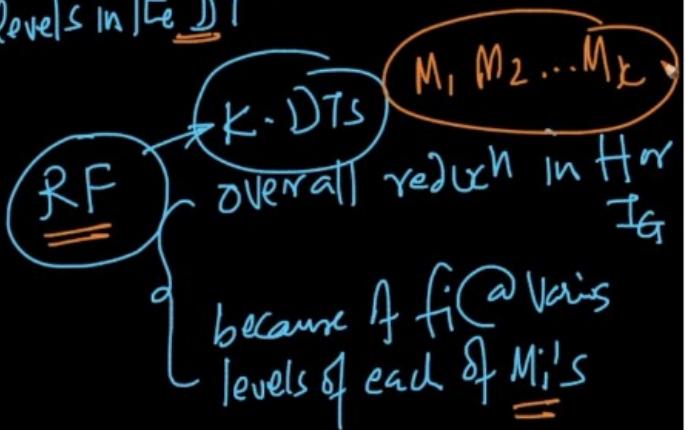
Feature importance: The overall reduction in entropy because of this feature at various levels in DT of all the base models.

Then the feature is said to be more important.

② Feature Imp:



overall reduction in Entropy or I_G
because of this feature @ various
levels in $[k] DT$



Boosting Intuition:

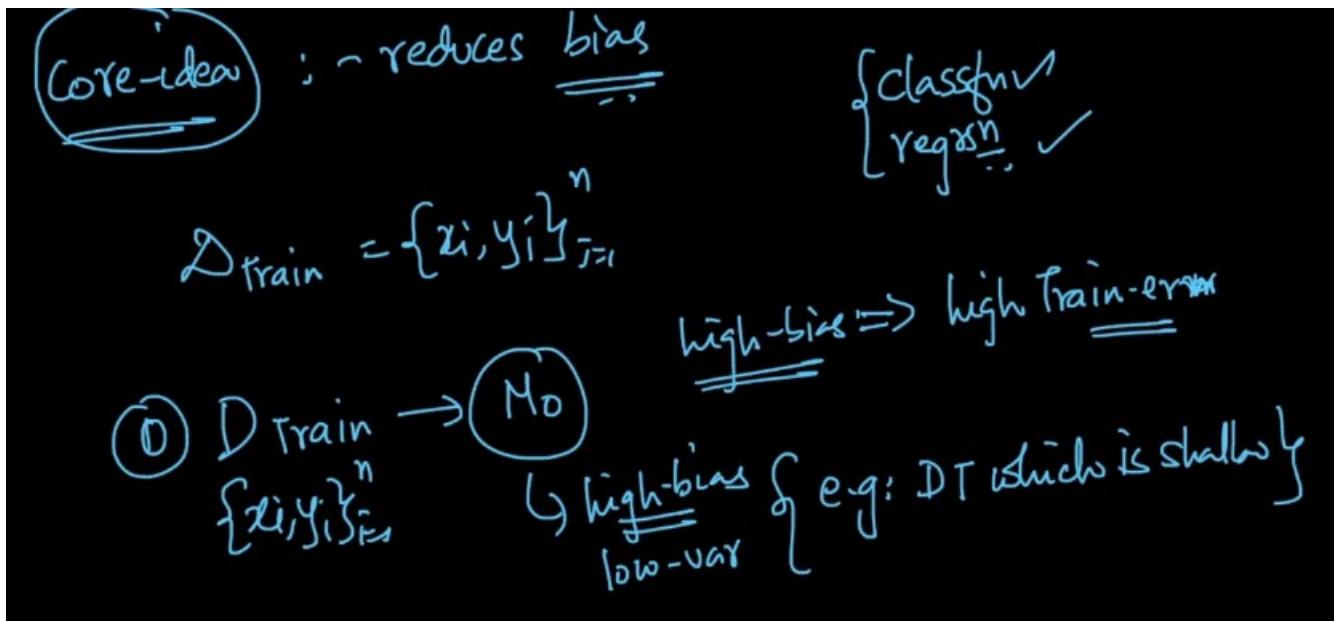
Instead of making high variance and low bias models in bagging, we use low variance and high bias models in Boosting with additive combine, that will reduce bias while keep out variance low.

✓ Boosting: intuition

Bagging: - high Var, low bias base-models + randomization + aggregation (CS, RS)

Boosting: low-var high bias + additively combine { reduce bias while keep our var low .. } GB

High bias can be viewed as high training error.



Step - 1:

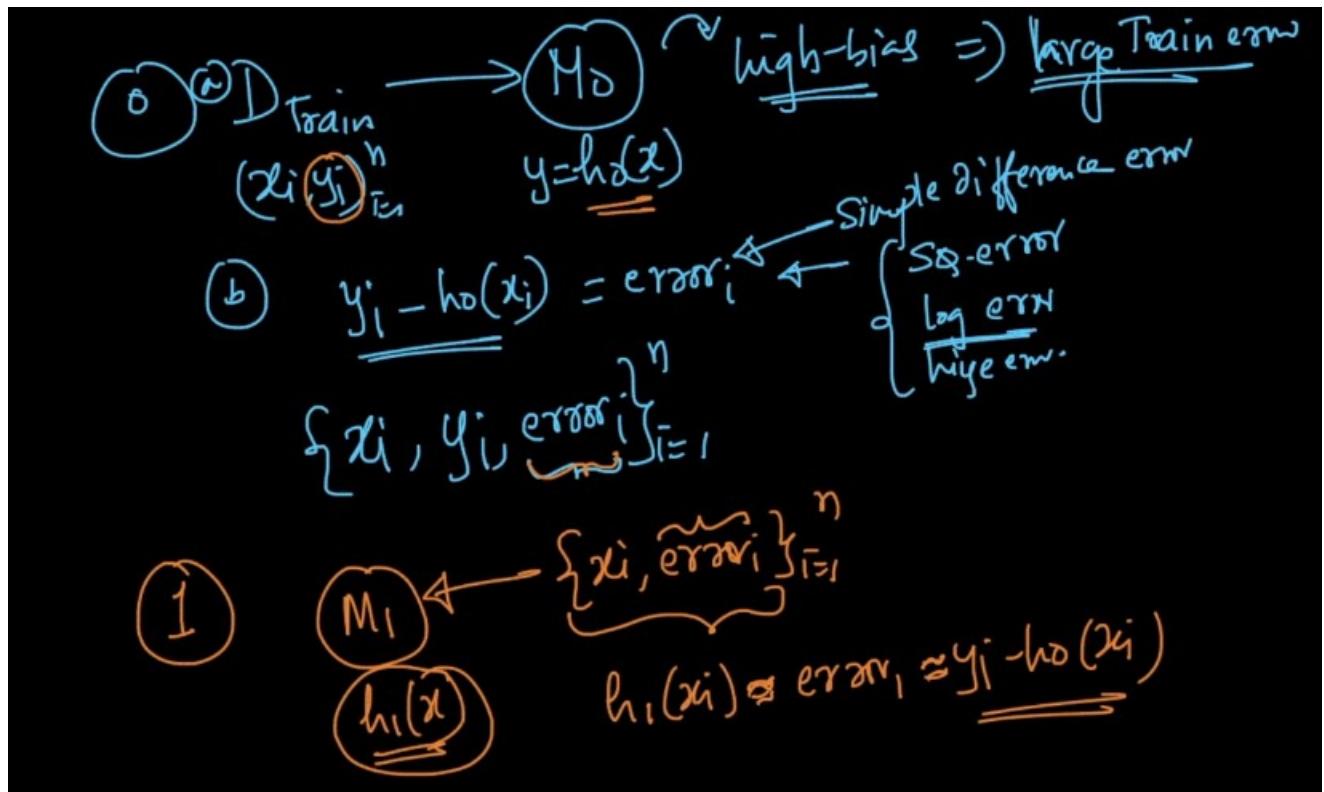
Train the model on all the data points. Make the model $h(x)$.

Step - 2:

Get the error on the points. $y_i - h(x_i)$. This is a simple error.

Step - 3:

Train the model on the error that we got from the previous model.



The models are made by additively combining the models.

$F_1(x)$ = model at end of Stage 1

$$F_1(x) = \underbrace{\alpha_0 h_0(x)}_{\text{constant}} + \underbrace{\alpha_1 h_1(x)}_{\text{base model}} : \text{- weighted sum of 2-base models}$$

② $\{x_i, err_i\} \rightarrow M_2$

$$y_i - F_1(x_i)$$

$$F_2(x) = \underbrace{\alpha_0 h_0(x)}_{\text{constant}} + \underbrace{\alpha_2 h_2(x)}_{\text{base model}} + \underbrace{\alpha_1 h_1(x)}_{\text{base model}}$$

Generalized notation:

end of stage k : -

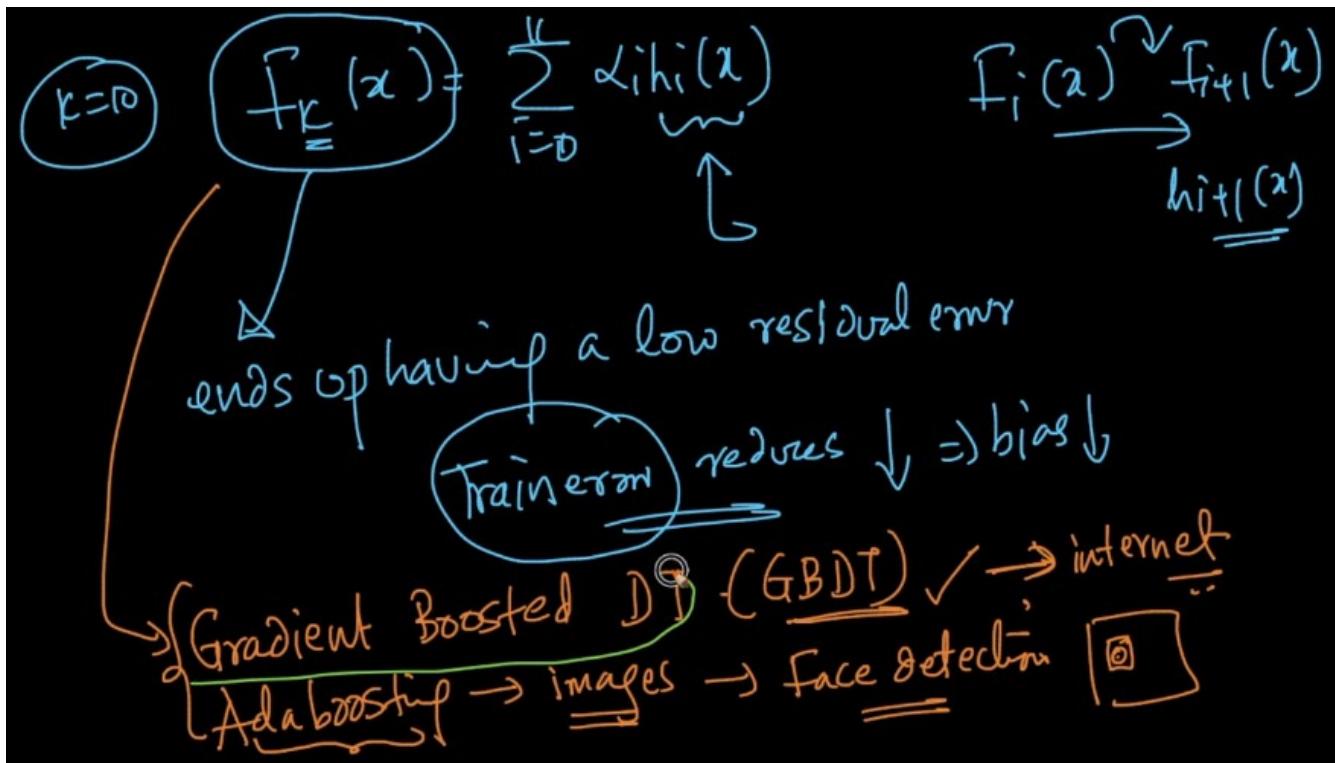
$$F_k(x) = \sum_{i=0}^k \underbrace{\alpha_i h_i(x)}_{\text{additive weighted model}} \quad \begin{array}{l} \text{trained to fit the} \\ \text{residual error @ end of} \\ \text{the prev stage} \end{array}$$

$h_i(x) \leftarrow \{x_i, err_i\}$

$y_i - F_{i-1}(x)$

$\uparrow \text{residual error @ end of stage } (i-1)$

The every model that is being constructed is to reduce the high error/bias. At every level the bias/error is reduced.



Residuals, loss functions and gradients:

Residuals:

$\overbrace{\text{Residuals, Loss-fns \& Gradients}} \rightarrow \overbrace{\begin{cases} \text{Gradient} \\ \text{Boosting} \end{cases}}$

$$f_r(x) = \sum_{i=0}^r \alpha_i h_i(x)$$

$M_{r+1} \leftarrow \{x_i, err_i\}$

$err_i = y_i - f_r(x)$
 Residual
 @ end of stage K

Here the goal is to find the α values and number of base learners. How to model the base learners.

Loss minimization:

For one base learner.

Loss-minimization:

- Logistic-loss \leftrightarrow classfn
- L2 reg \leftrightarrow Sq. loss \leftrightarrow easier
- SVM \leftrightarrow hinge-loss

regression:

$$L(y_i, \tilde{f}_k(x_i)) = (y_i - \tilde{f}_k(x_i))^2$$

$\left\{ \begin{array}{l} \text{let} \\ \tilde{f}_k(x_i) = z_i \end{array} \right.$

Sq. loss

$$\frac{\partial L}{\partial \tilde{f}_k(x_i)} = \frac{\partial L}{\partial z_i} = \frac{\partial}{\partial z_i} (y_i - z_i)^2$$

$$= (-1) * 2 * (y_i - z_i)$$

$$= -2(y_i - z_i)$$

At the end of the 'K' th base model.

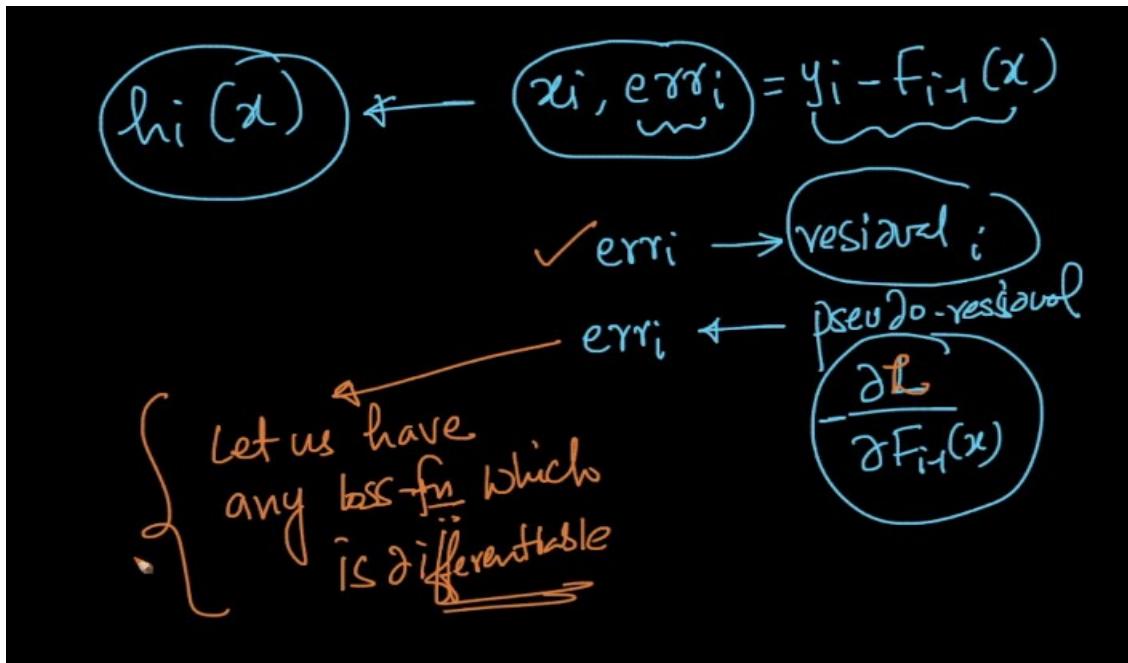
The error is called the residual error. The negative gradient can be taught like it is proportional to the residual error.

$$-\frac{\partial \tilde{L}}{\partial F_k(x_i)} = \boxed{\times} \underbrace{(y_i - \tilde{F}_k(x_i))}_{\text{residual}} \quad \begin{matrix} \text{error} \\ \text{residual} \end{matrix}$$

neg derivative

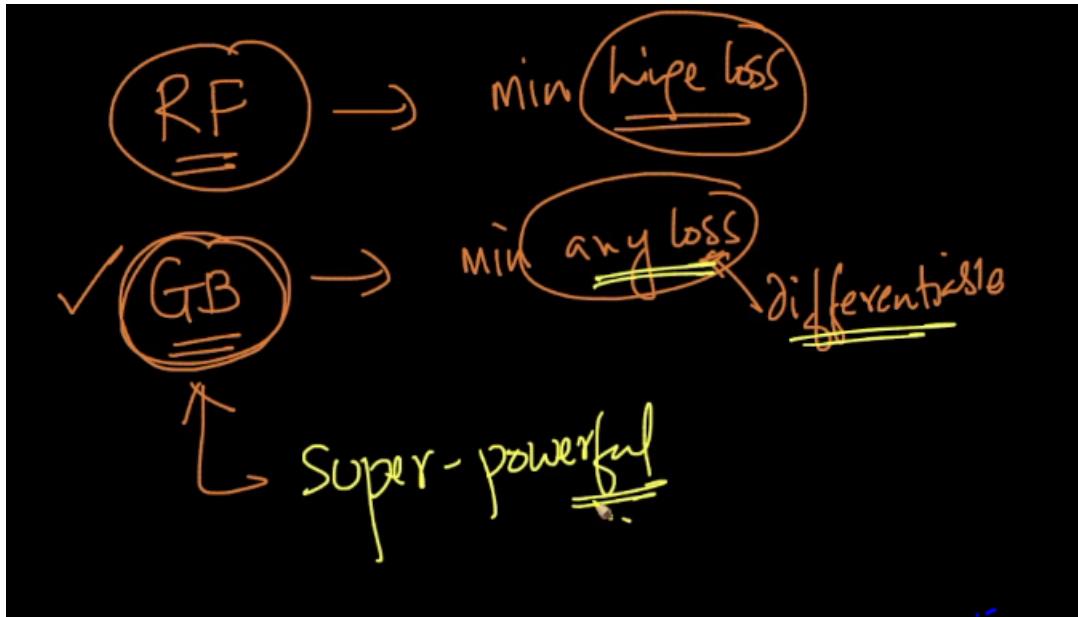
* neg. gradient \approx residual

This is often called as pseudo residuals.

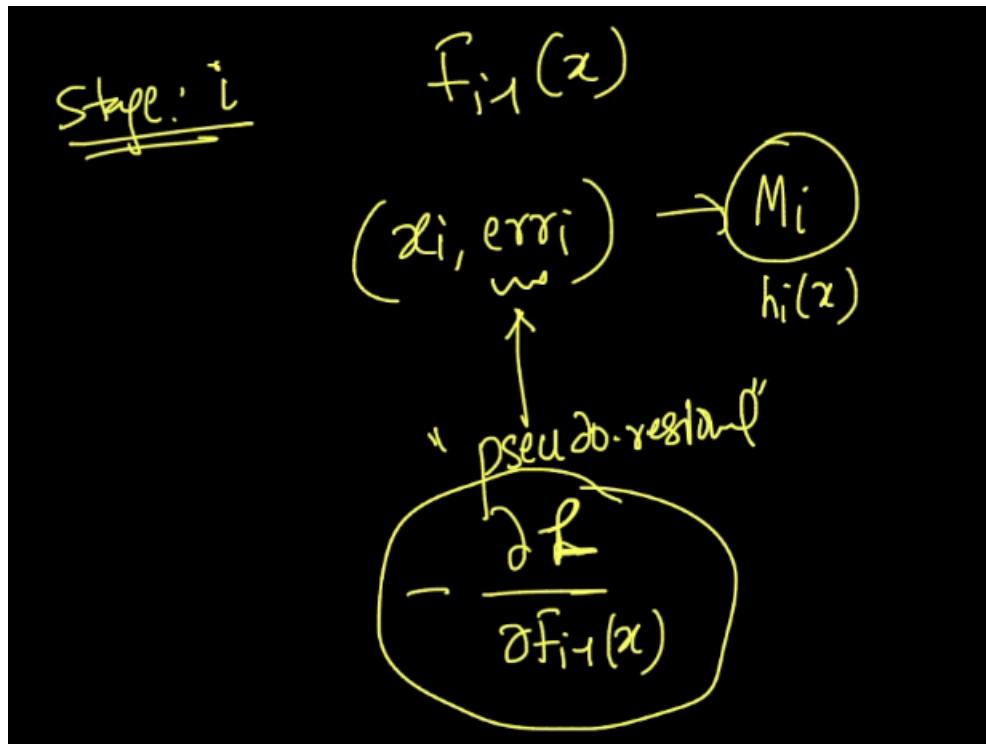


Replacing the residuals with the pseudo residuals, will help in replacing with the any loss function which is differential.

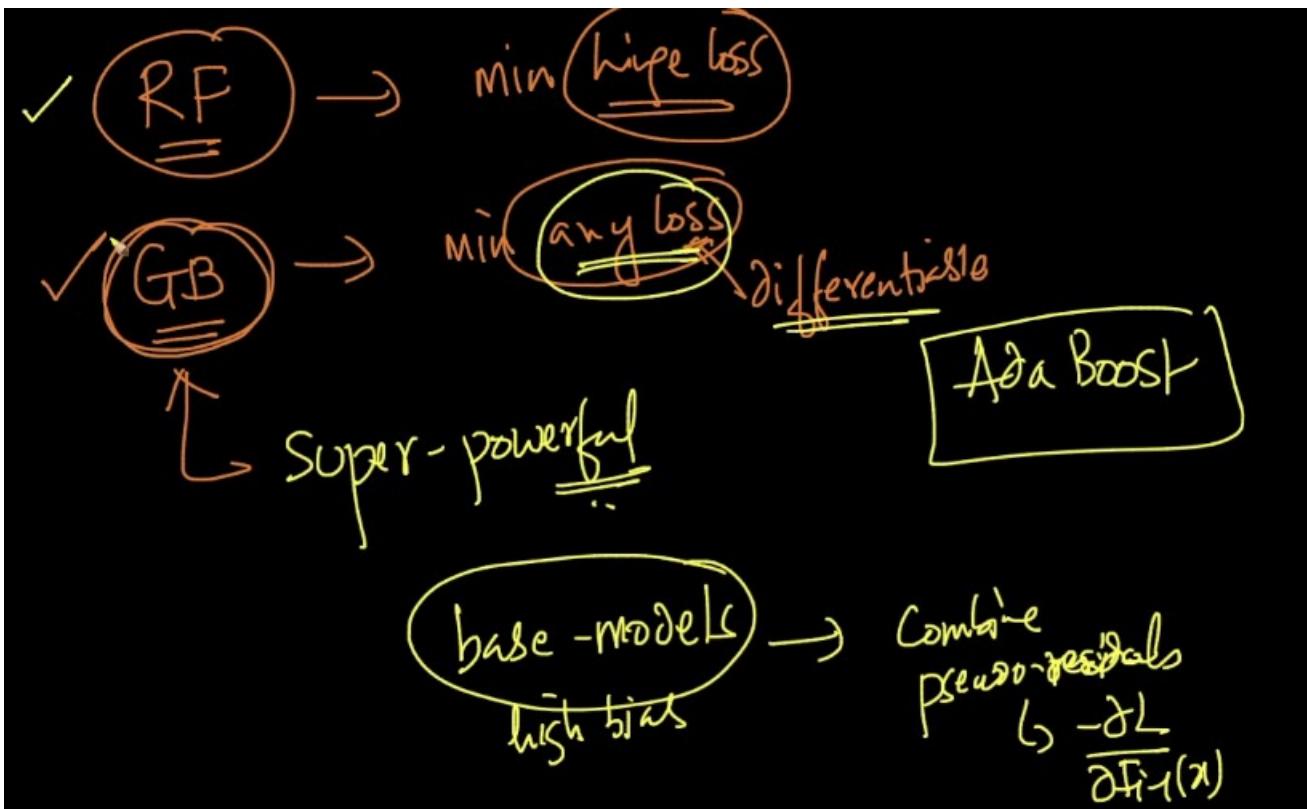
As long as the loss function is differential we can replace with any other loss function. This is the reason gradient boosting is so popular.



At any stage 'i' we can take the pseudo residual of the actual error got from the previous level.



We have the other algorithm called Ada-boost.



Gradient Boosting:

Input: training set $\{(x_i, y_i)\}_{i=1}^n$, a differentiable loss function $L(y, F(x))$, number of iterations M .

Algorithm:

1. Initialize model with a constant value:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$

2. For $m = 1$ to M :

1. Compute so-called *pseudo-residuals*:

$$r_{im} = - \left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$

2. Fit a base learner (e.g. tree) $h_m(x)$ to pseudo-residuals, i.e. train it using the training set $\{(x_i, r_{im})\}_{i=1}^n$.

3. Compute multiplier γ_m by solving the following one-dimensional optimization problem:

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$

4. Update the model:

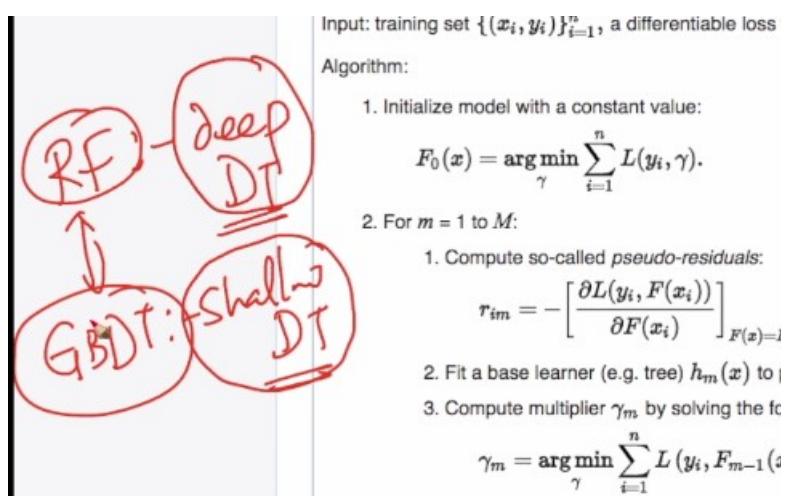
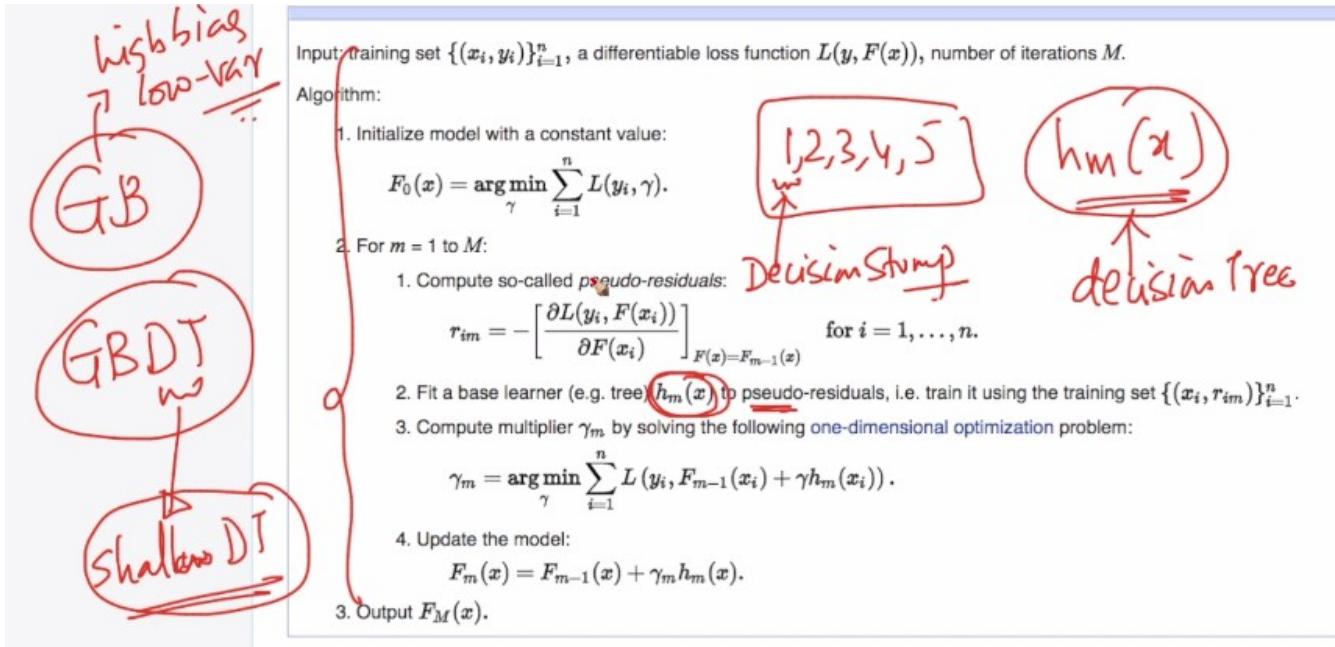
$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

3. Output $F_M(x)$.

$$F_M(x) = \underbrace{\gamma_1 h_1(x)}_{\gamma_1 h_1(x)} + \underbrace{\gamma_2 h_2(x)}_{\gamma_2 h_2(x)} + \dots + \underbrace{\gamma_m h_m(x)}_{\gamma_m h_m(x)}$$

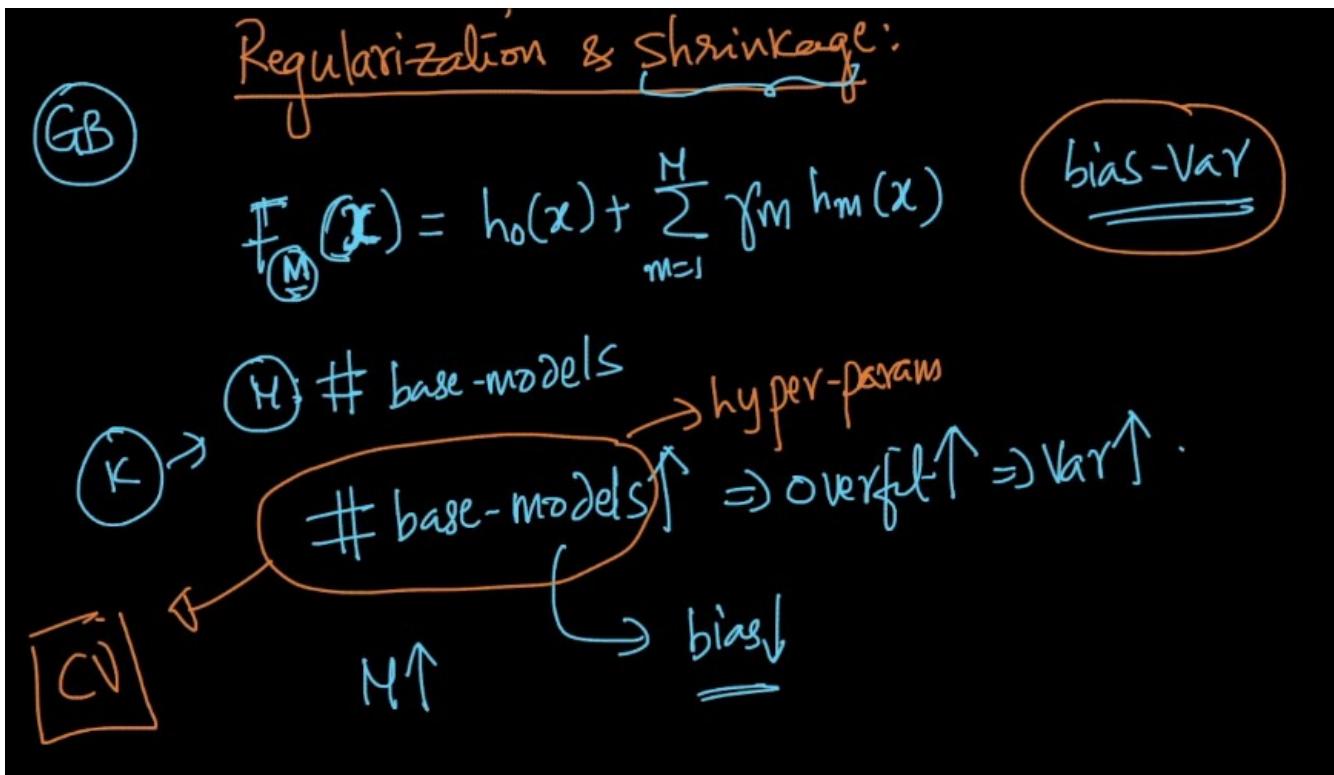
mth iter.

Gradient boosting is the general idea. The each of the model must be high bias and low variance model. This can occur in case of decision trees with shallow depths.



Regularization and shrinkage:

As the number of base models increase then we would over fit and variance can also increase, bias can decrease.



The number of base models are the hyper parameter, that can be chosen by cross – validation.

Here the shrinkage can be tuned by using the variable 'V', which is called the learning rate. Lies between 0 and 1.

Regularization [edit]

Fitting the training set too closely can lead to degradation of the model's generalization ability. Several so-called regularization techniques reduce this overfitting effect by constraining the fitting procedure.

One natural regularization parameter is the number of gradient boosting iterations M (i.e. the number of trees in the model when the base learner is a decision tree). Increasing M reduces the error on training set, but setting it too high may lead to overfitting. An optimal value of M is often selected by monitoring prediction error on a separate validation data set. Besides controlling M , several other regularization techniques are used.

Shrinkage [edit]

An important part of gradient boosting method is regularization by shrinkage which consists in modifying the update rule as follows:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x)$$

$$F_m(x) = F_{m-1}(x) + \nu \cdot \gamma_m h_m(x), \quad 0 < \nu \leq 1, \quad \nu \approx 0.1$$

where parameter ν is called the "learning rate".

Empirically it has been found that using small learning rates (such as $\nu < 0.1$) yields dramatic improvements in model's generalization ability over gradient boosting without shrinking ($\nu = 1$).^[7] However, it comes at the price of increasing computational time both during training and querying: lower learning rate requires more iterations.

Stochastic gradient boosting [edit]

Soon after the introduction of gradient boosting Friedman proposed a minor modification to the algorithm, motivated by Breiman's bagging method.^[3] Specifically, he proposed that at each iteration of the algorithm, a base learner should be fit on a random subset of the training data (the bootstrap sample).

Here 'V' is also a hyper parameter.

If 'V' is very small then the chances of over fitting also reduces and variance also decreases.

If 'V' increases then the chances of overfitting also increases.

M is the number of base models.

Regularization [edit]

Fitting the training set too closely can lead to degradation of the model's generalization ability. Several so-called regularization techniques reduce this overfitting effect by constraining the fitting procedure.

One natural regularization parameter is the number of gradient boosting iterations M (i.e. the number of trees in the model when the base learner is a decision tree). Increasing M reduces the error on training set, but setting it too high may lead to overfitting. An optimal value of M is often selected by monitoring prediction error on a separate validation data set. Besides controlling M , several other regularization techniques are used.

\downarrow
Overfitting
 \downarrow
Want

Shrinkage [edit]

An important part of gradient boosting method is regularization by shrinkage which consists in modifying the update rule as follows:

$$F_m(x) = F_{m-1}(x) + \nu \underline{\gamma_m} h_m(x), \quad 0 < \nu \leq 1,$$

where parameter ν is called the "learning rate".

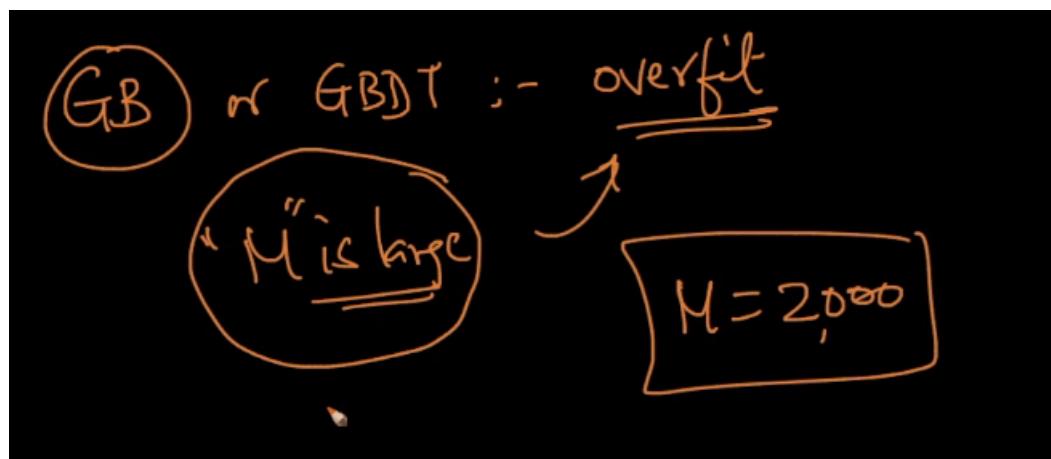
$M = \# \text{base-models}$
 $\uparrow \Rightarrow \text{overfitting} \uparrow$

Empirically it has been found that using small learning rates (such as $\nu < 0.1$) yields dramatic improvements in model's generalization ability over gradient boosting without shrinking ($\nu = 1$).^[7] However, it comes at the price of increasing computational time both during training and querying: lower learning rate requires more iterations.

Stochastic gradient boosting [edit]

The ideal method is to find the best parameters using grid search.

It is easy to over-fit for the gradient boosted trees.



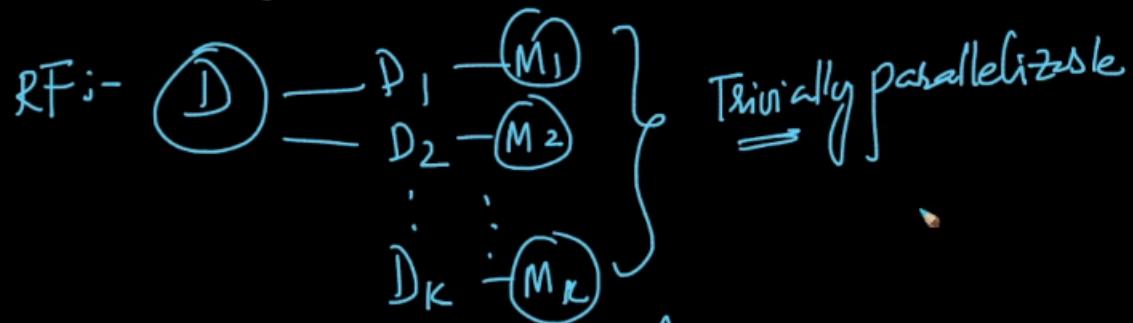
The train and test time complexity:

GBDT is not easy to parallelize.

Train & Run time complexity

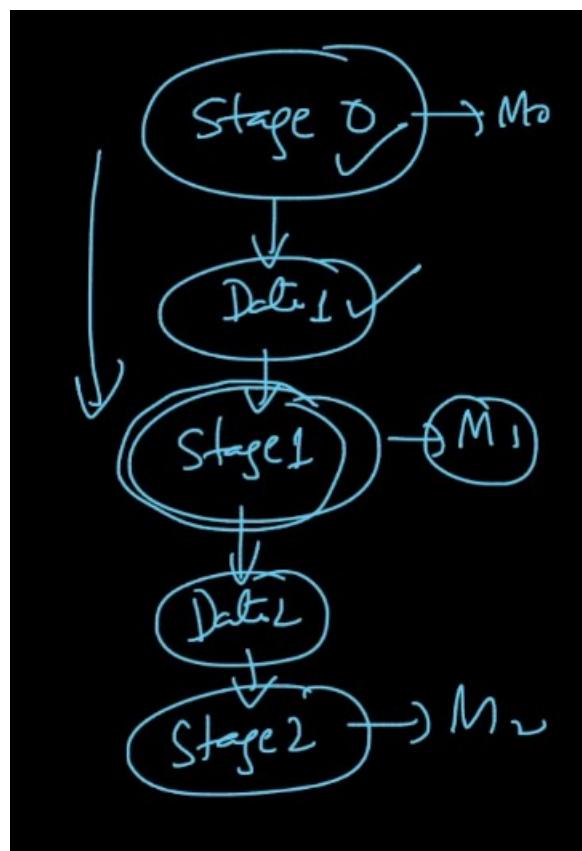


Train: $O(n \lg n d * M)$ M : #base-learners



GBDT:- not easy to parallelize:-

In GBDT all the models are related to each other. The data is also produced after each of the stage.



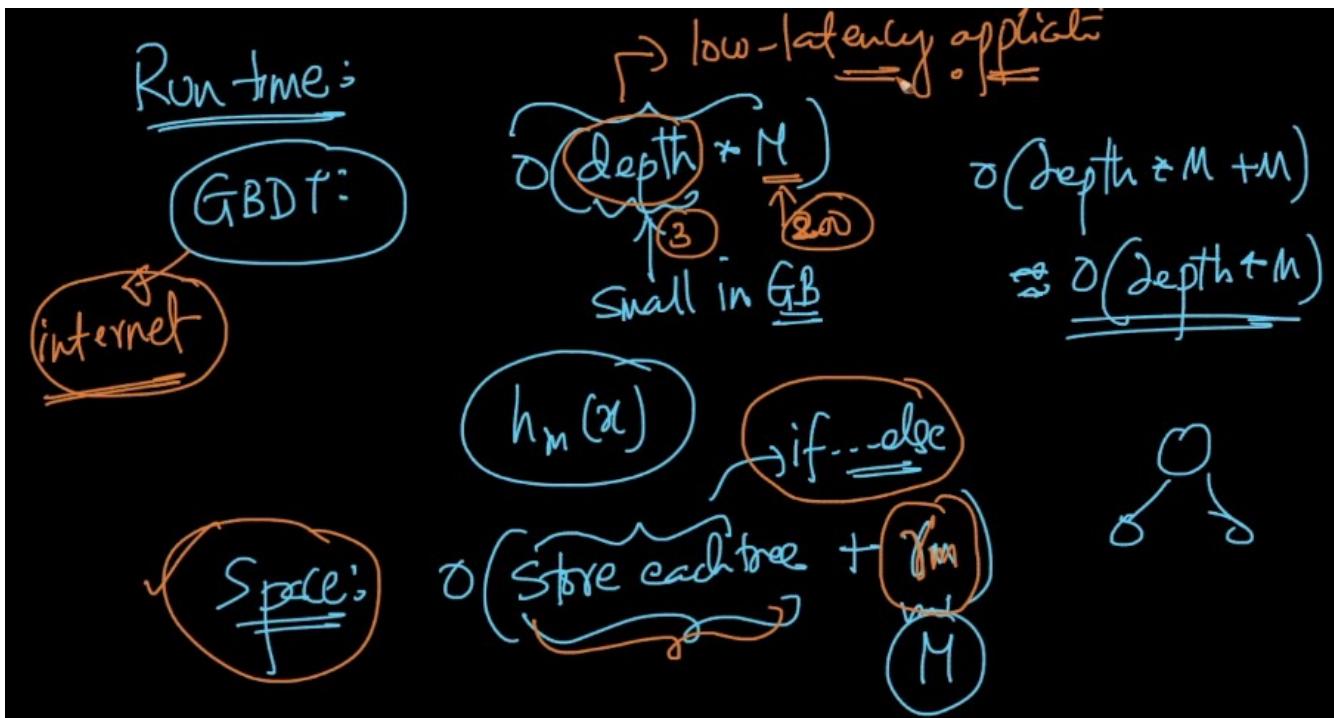
Hence GBDT is not trivial to parallelize.

GBDT takes more time than random forest, even the time complexity remains same.



Run time:

At run time the depth is very low and one must keep track of all the coefficient values.



Gradient Boosting Classifier(Sklearn):

Scikit-learn Documentation

3.2.4.3.5. `sklearn.ensemble.GradientBoostingClassifier`

```
class sklearn.ensemble.GradientBoostingClassifier(loss='deviance', learning_rate=0.1, n_estimators=100,
subsample=1.0, criterion='friedman_mse', min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0,
max_depth=3, min_impurity_decrease=0.0, min_impurity_split=None, init=None, random_state=None,
max_features=None, verbose=0, max_leaf_nodes=None, warm_start=False, presort='auto')
```

Gradient Boosting for classification.

GB builds an additive model in a forward stage-wise fashion; it allows for the optimization of arbitrary differentiable loss functions. In each stage `n_classes_` regression trees are fit on the negative gradient of the binomial or multinomial deviance loss function. Binary classification is a special case where only a single regression tree is induced.

Read more in the User Guide.

Parameters: `loss` : {‘deviance’, ‘exponential’}, optional (default=‘deviance’)

loss function to be optimized. ‘deviance’ refers to deviance (= logistic regression) for

Annotations:

- A red circle with 'GB' is drawn over the class name `GradientBoostingClassifier`.
- A red circle with 'shallow' is drawn over the text 'shallow'.
- A red checkmark is drawn next to the code snippet.
- A red circle with 'M' is drawn over the parameter `n_estimators`.
- A red circle with 'Y' is drawn over the parameter `loss`.
- A red circle with 'N' is drawn over the parameter `criterion`.

The sub sample can be taught of the row_sampling parameter.

The minimum number of samples required to be at a leaf node:

- If int, then consider min_samples_leaf as the minimum number.
- If float, then min_samples_leaf is a percentage and ceil(min_samples_leaf * n_samples) are the minimum number of samples for each node.

Changed in version 0.18: Added float values for percentages.

min_weight_fraction_leaf : float, optional (default=0.)

The minimum weighted fraction of the sum total of weights (of all the input samples) required to be at a leaf node. Samples have equal weight when sample_weight is not provided.

subsample : float, optional (default=1.0)

The fraction of samples to be used for fitting the individual base learners. If smaller than 1.0 this results in Stochastic Gradient Boosting. subsample interacts with the parameter n_estimators. Choosing subsample < 1.0 leads to a reduction of variance and an increase in bias.

max_features : int, float, string or None, optional (default=None)

The number of features to consider when looking for the best split:

- If int, then consider max_features features at each split.
- If float, then max_features is a percentage and int(max_features * n_features) fea-

The value loss= , is the measure of type of loss function to be used.
If loss = ‘AdaBoost’, this algorithm is equal to another algorithm called ada – boost.

scikit-learn.org/stable/modules/generated/sklearn.ensemble.GradientBoostingClassifier.html

3.2.4.3.5. `sklearn.ensemble.GradientBoostingClassifier`

scikit-learn v0.19.1
Other versions

Please cite us if you use the software.

3.2.4.3.5.
`sklearn.ensemble.GradientBoostingClassifier`

`class sklearn.ensemble.GradientBoostingClassifier(loss='deviance', learning_rate=0.1, n_estimators=100, subsample=1.0, criterion='friedman_mse', min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_depth=3, min_impurity_decrease=0.0, min_impurity_split=None, init=None, random_state=None, max_features=None, verbose=0, max_leaf_nodes=None, warm_start=False, presort='auto')` [source]

Gradient Boosting for classification.

GB builds an additive model in a forward stage-wise fashion; it allows for the optimization of arbitrary differentiable loss functions. In each stage n_classes_ regression trees are fit on the negative gradient of the binomial or multinomial deviance loss function. Binary classification is a special case where only a single regression tree is induced.

Read more in the User Guide.

Parameters:

- loss** : {‘deviance’ | ‘exponential’}, optional (default=‘deviance’)
- loss function to be optimized. ‘deviance’ refers to deviance (= logistic regression) for classification with probabilistic outputs. For loss ‘exponential’ gradient boosting recovers the AdaBoost algorithm.
- learning_rate** : float, optional (default=0.1)
- learning rate shrinks the contribution of each tree by learning_rate. There is a trade-off between learning_rate and n_estimators.

The sklearn implementation is not so efficient.
It is very slow.

Xgboost is the efficient way implementation.(well optimized).

XGBoost Get Started Tutorials How To Packages Knobs Search

Scikit-Learn API

Scikit-Learn Wrapper interface for XGBoost.

```
class xgboost.XGBRegressor(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True, objective='reg:linear', booster='gbtree', n_jobs=1, nthread=None, gamma=0, min_child_weight=1, max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_alpha=0, reg_lambda=1, scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None, **kwargs)
```

Bases: `xgboost.sklearn.XGBModel, object`

Implementation of the scikit-learn API for XGBoost regression.

Parameters

- max_depth : int**
Maximum tree depth for base learners.
- learning_rate : float**
Boosting learning rate (xgb's "eta")
- n_estimators : int**

In XGBoost we use the GBDT + Row sampling + column sampling. XGBoost data sets perform well than the sklearn implementation due to several code optimizations.

XGBoost Get Started Tutorials How To Packages Knobs Search

type:

Scikit-Learn API

Scikit-Learn Wrapper interface for XGBoost.

```
class xgboost.XGBRegressor(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True, objective='reg:linear', booster='gbtree', n_jobs=1, nthread=None, gamma=0, min_child_weight=1, max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_alpha=0, reg_lambda=1, scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None, **kwargs)
```

Bases: `xgboost.sklearn.XGBModel, object`

Implementation of the scikit-learn API for XGBoost regression.

Parameters

- max_depth : int**
Maximum tree depth for base learners.
- learning_rate : float**
Boosting learning rate (xgb's "eta")

There is the choice of choosing the features at the level of tree and at the level each level in the tree.

There is the similar representation in case of linear regression and GB. We can look <gamma's> as the weights.

XGBoost Get Started Tutorials How To Packages - Knobs Search

colsample_bylevel : float
Subsample ratio of columns for each split, in each level.

reg_alpha : float (xgb's alpha)
L1 regularization term on weights

reg_lambda : float (xgb's lambda)
L2 regularization term on weights

scale_pos_weight : float
Balancing of positive and negative weights.

base_score:
The initial prediction score of all instances, global bias.

seed : int
Random number seed. (Deprecated, please use random_state)

random_state : int
Random number seed. (replaces seed)

missing : float, optional

Sparse GBDT

$$f_m \rightarrow \text{weights}$$

$$w_0 + w_1 x_{11} + w_2 x_{12} + w_3 x_{13} + \dots$$

$$h_0(x) + h_1(x) + h_2(x) + h_3(x)$$

$L_1 \cdot reg := loss + \underline{\lambda_{reg}}$

$$\alpha L_1 + \lambda b_2$$

Ada-boost geometric intuition:
It is mainly used in case of image data for face detection.

A toy example from Schapire's tutorial

Consider this toy dataset in 2D and assume that our weak classifiers are decision stumps (vertical or horizontal half-planes):

D_1

GBDT

Boosting! - Ada-boost

The first round:

Image

face detection

h_1

$\epsilon_1 = 0.30$
 $\alpha_1 = 0.42$

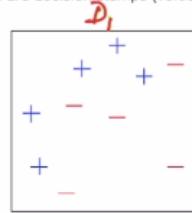
The second round:

h_2

h_3

A toy example from Schapire's tutorial

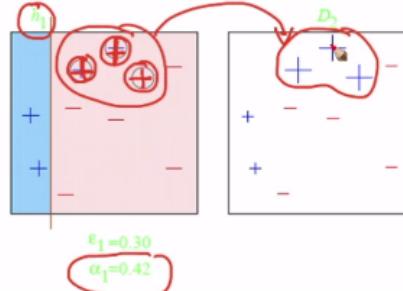
Consider this toy dataset in 2D and assume that our weak classifiers are decision stumps (vertical or horizontal half-planes):



DT. with depth = 1
Decision Stump

The first round:

$d_1 h_1(x)$



weighted models
↓
upsample

The second round:

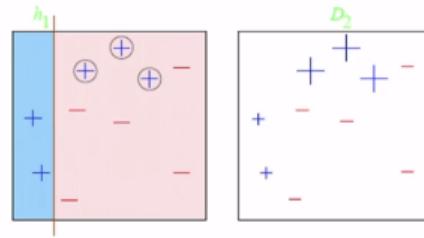


The error points in the data are up sampled in the next level.

The weights are increased exponentially in adaptive boosting.

Even of the errors we also get the alpha2 as the weight for the error points.

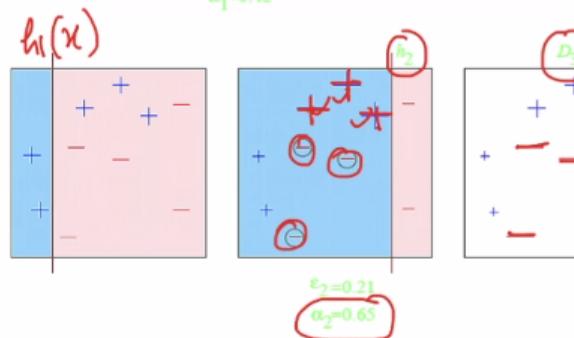
$\alpha_1 h_1$



~~exponentially~~

The second round:

$\alpha_2 h_2$



The third round:

In this way we implement more models with there weights.

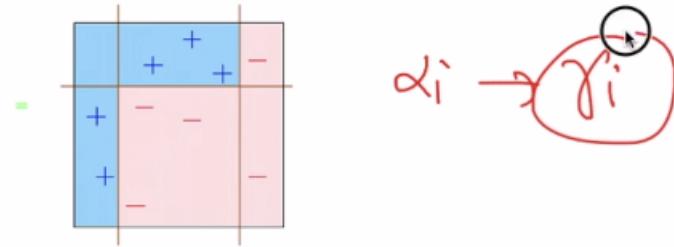


The final classifier

$$f_3(x) = \alpha_1 h_1 + \alpha_2 h_2 + \alpha_3 h_3$$

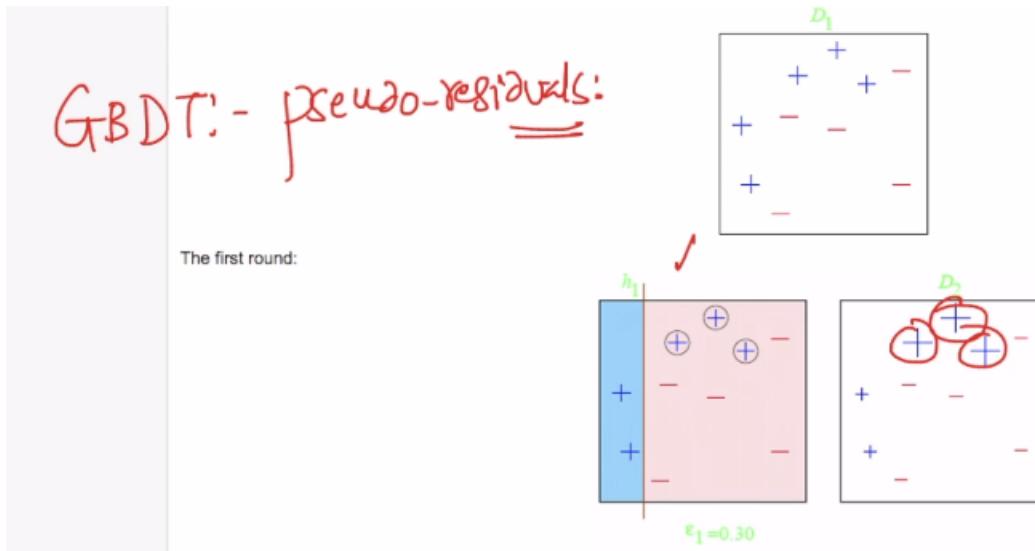
$$H_{\text{final}} = \text{sign} \left(\alpha_1 \cdot 0.42 + \alpha_2 \cdot 0.65 + \alpha_3 \cdot 0.92 \right)$$

The diagram shows the final classifier H_{final} as a sum of weighted base functions h_1 , h_2 , and h_3 . The weights are $\alpha_1 = 0.42$, $\alpha_2 = 0.65$, and $\alpha_3 = 0.92$. The final result is the sign of the sum of these weighted functions.



In case of GBDT we use the pseudo – residuals.

In case if ADABOOST, we reduce the error-nous points.



The weights are increased exponentially and the miss-classified points are up sampled at each level.

Ada-boost is applied very less in real – world. It is mainly used for face detection.

Stacking models:

Mlxtend classifier for implementing stacking models.

In stacking the data is trained on several classification algorithms, the algos must be very different from each other. These are called the base models.(first – level classifiers). The more different the models are the better it works.

Then again construct the data set on the predictions made by the base models for every point in the data set.

Then on the new data set a classifier is trained to make the final output of the query point.

Algorithm 19.7 Stacking

Input: Training data $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^m$ ($\mathbf{x}_i \in \mathbb{R}^n$, $y_i \in \mathcal{Y}$)

Output: An ensemble classifier H

```
1: Step 1: Learn first-level classifiers
2: for  $t \leftarrow 1$  to  $T$  do
3:   Learn a base classifier  $h_t$  based on  $\mathcal{D}$ 
4: end for
5: Step 2: Construct new data sets from  $\mathcal{D}$ 
6: for  $i \leftarrow 1$  to  $m$  do
7:   Construct a new data set that contains  $\{\mathbf{x}'_i, y_i\}$ , where  $\mathbf{x}'_i = \{h_1(\mathbf{x}_i), h_2(\mathbf{x}_i), \dots, h_T(\mathbf{x}_i)\}$ 
8: end for
9: Step 3: Learn a second-level classifier
10: Learn a new classifier  $h'$  based on the newly constructed data set
11: return  $H(\mathbf{x}) = h'(h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_T(\mathbf{x}))$ 
```

The last classifier is called the meta classifier.

Cascading models:

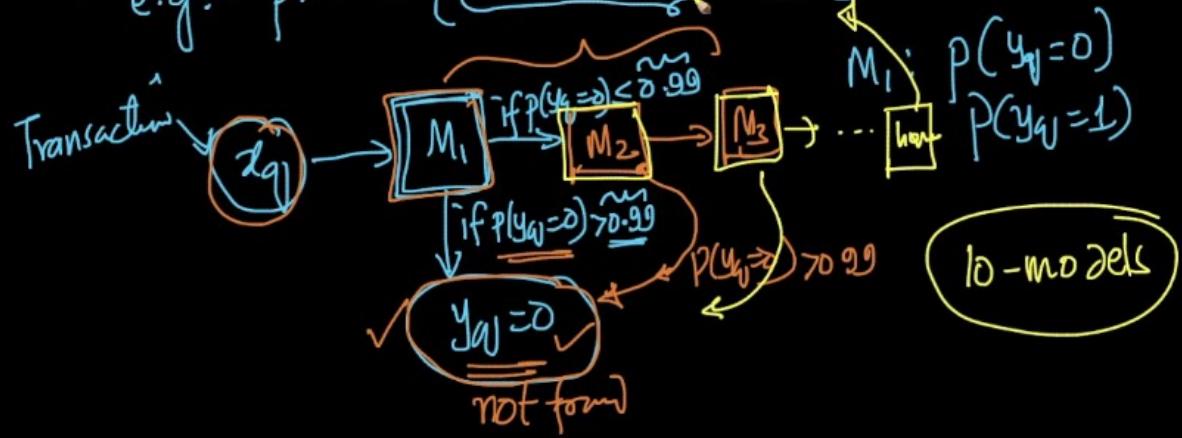
Example for understanding the cascading models. Every transaction is a vector, every fraud transaction can be done from various locations.

Every model can give the class probabilities. (fraud and not – fraud). If probability of the point to be one class can be large and other class is low.

At every stage we train an another classifier to check the certainty of the model.

"Cascading" models:

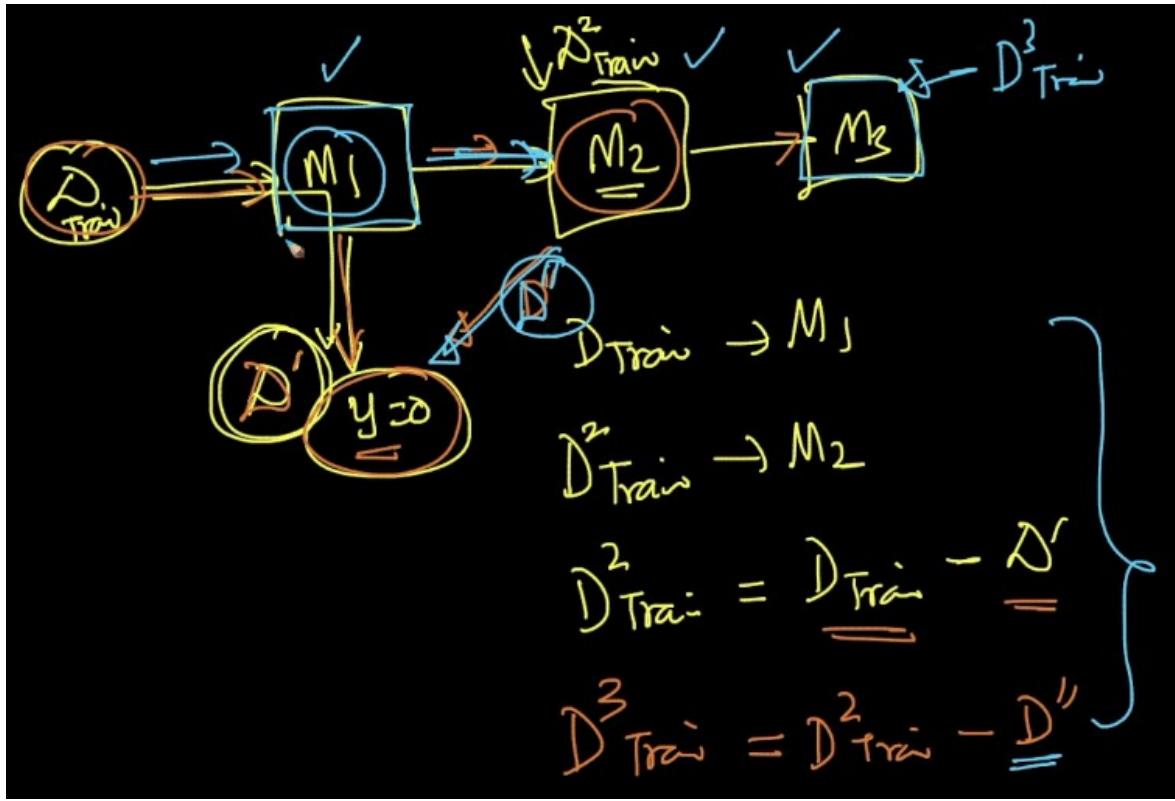
- Bagging, boosting, Stacking
- e.g.: predicts {credit card transaction} is fraudulent or not



They are typically used when the cost of making mistake is very high.

The points that don't give the accuracy that is being passed to the other model for achieving the accuracy and so on.. till there can be a specialist checking at the end.

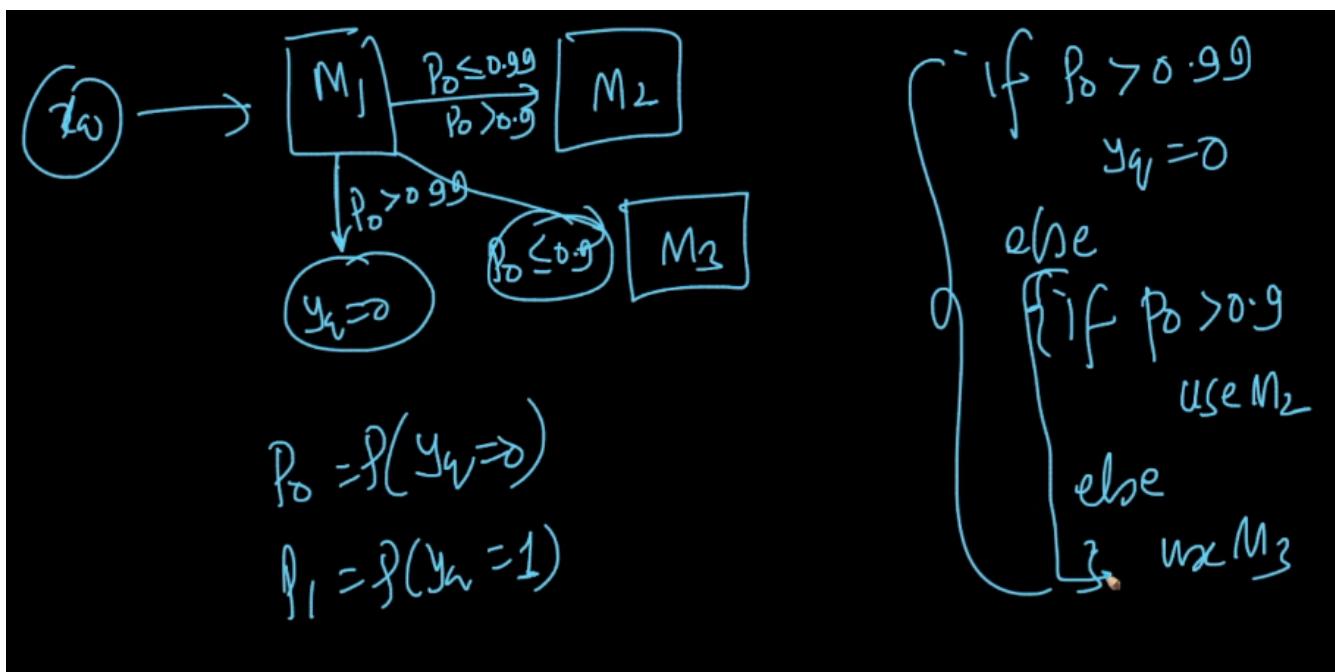
At every level the data points in the data set are reduced.



The cascade can be more than one type.

The down output of the model is considered to be the class of the model based on the probability of the model.(i.e, say the model outputs 0.99 as the probability score, then the model will confirm the class of the data point, the remaining points are further sent to the cascade model of the first model).

We could make the complex if – else block statements for predicting the models outcome based on the accuracy that is needed.



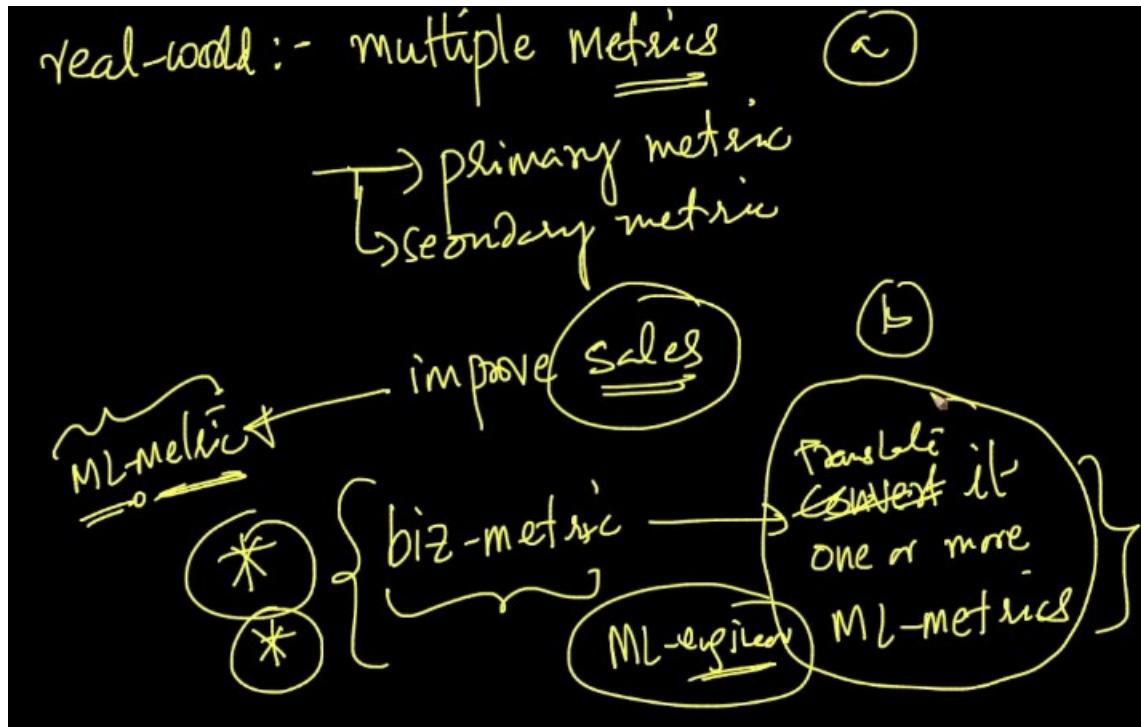
Cascades models can be extensively used in fraud detection, medicine applications for maintaining the accuracy of the model.

They can also be build up to 20 levels.

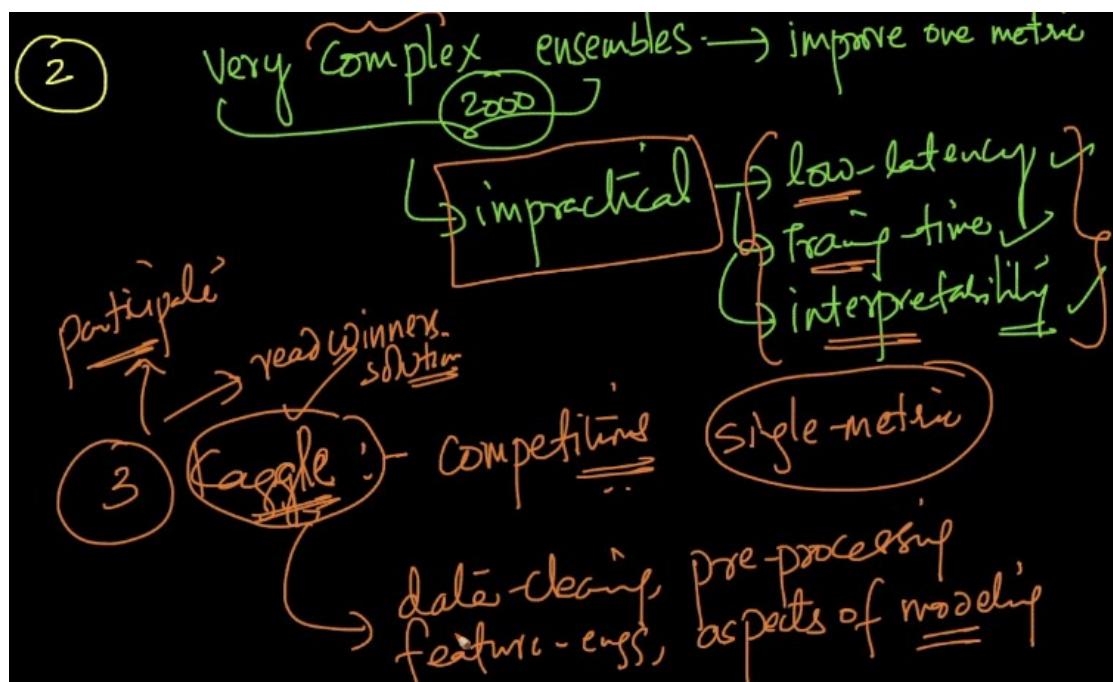
If one transaction is lost then there will be huge loss.

Kaggle vs. Real – world cases:

Kaggle can only take one metric as the measure. A real world there are multiple metrics that can be taken care of that favors the business problem. They take the business goal and translate the biz-problem to one – more machine learning metrics.



In case – studies we will see the translation of the biz – problem to a machine learning metrics. They build the very complex ensembles, this may be impractical.



Exercise:

