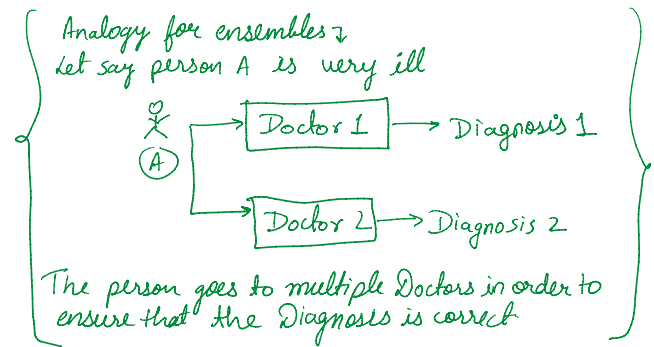
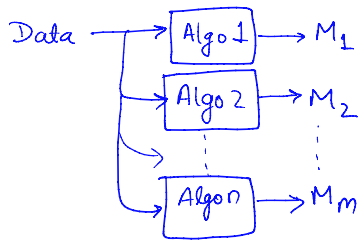


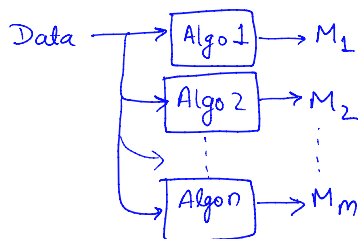
## Ensembles

means a group of objects  
↓  
models



## Ways to Build Ensembles

### ① Parallel Ensembles



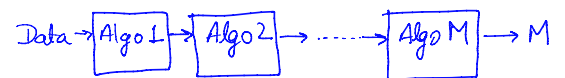
1.1 Voting Ensembles

1.2 Stacking (sklearn)

1.3 Bagging

Random Forest Algo. sklearn

### ② Sequential Ensembles



1.1 Cascading

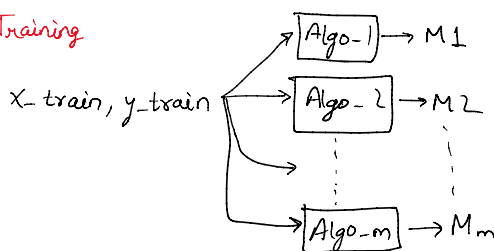
1.2 Boosting



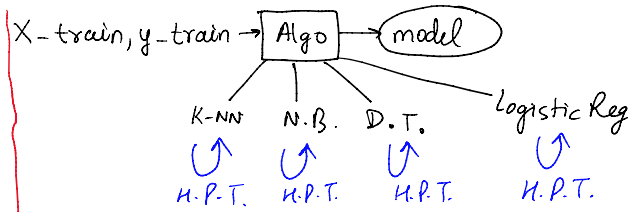
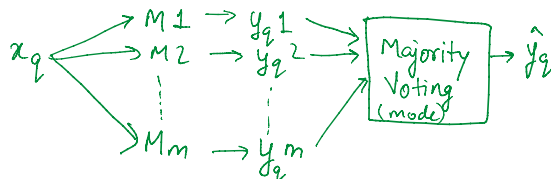
## Voting Ensemble

Classification Task →

Training



Prediction



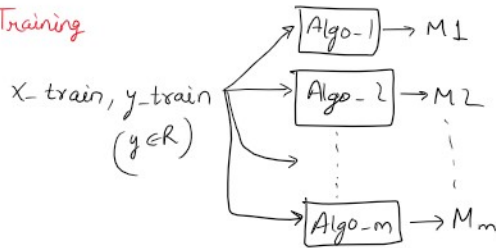
Prediction →



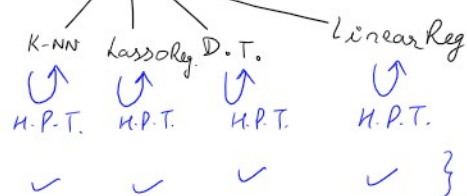
} Final model selection from the best is made on the basis of requirements.

## Regression Task

Training

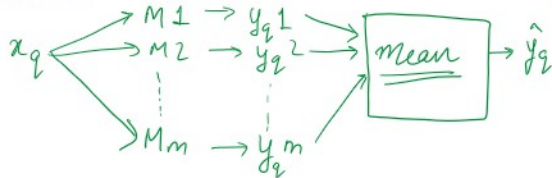


$X_{\text{train}}, y_{\text{train}} \rightarrow \text{Algo} \rightarrow \text{model}$



} Final model selection from the best is made on the basis of requirements.

Prediction

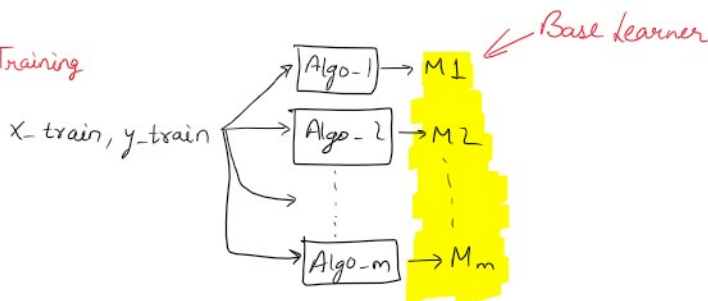


Prediction  $\rightarrow$

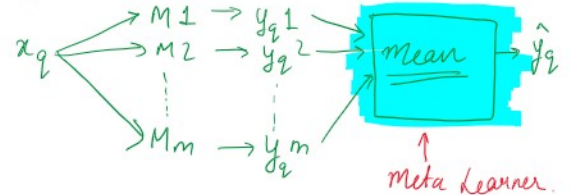


## Important Terminologies for Ensembles

Training



Prediction



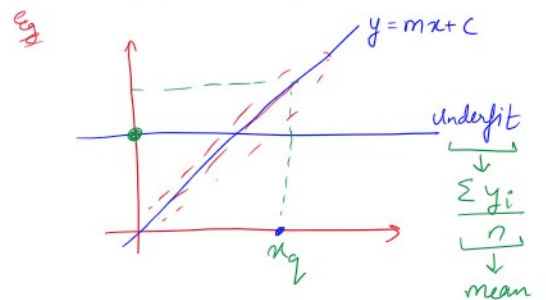
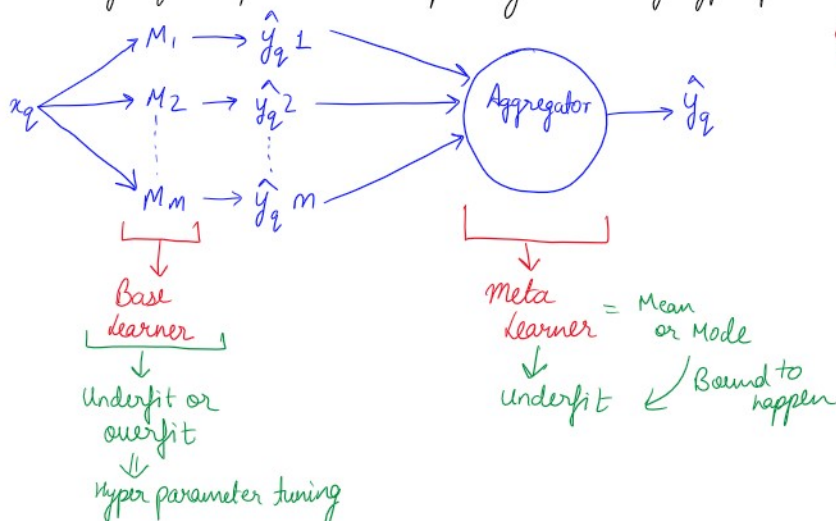
In Voting Ensembles

Classification { Base learner  $\rightarrow$  K-NN, D.T., N.B., logistic Reg etc.  
Meta learner  $\rightarrow$  Mode i.e. Majority Voting

Regression { Base learner  $\rightarrow$  K-NN, D.T., linear Reg., Lasso Reg., Theil-Sen Reg. etc.  
Meta learner  $\rightarrow$  Mean or Median

## Issue with Voting Ensemble

$\rightarrow$  It has very high computational complexity because of hyper parameter tuning for each algorithm in the Base learner.



Voting ensemble has the problem of  $\rightarrow$  overfitting & underfitting + Hyperparameter Tuning of multiple models which made it computationally very complex

As a 'solution' to the problems of Voting Ensembles  $\rightarrow$  we go to a stronger ensemble method

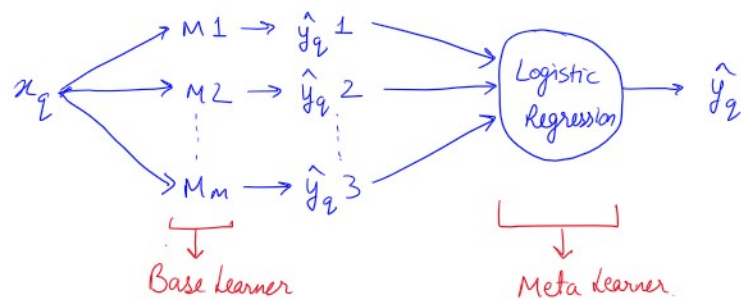
Stacking

Stacking

Classification  $\rightarrow$

Training  $\rightarrow$  Same as Voting Ensembles (Same in Regression as Well)

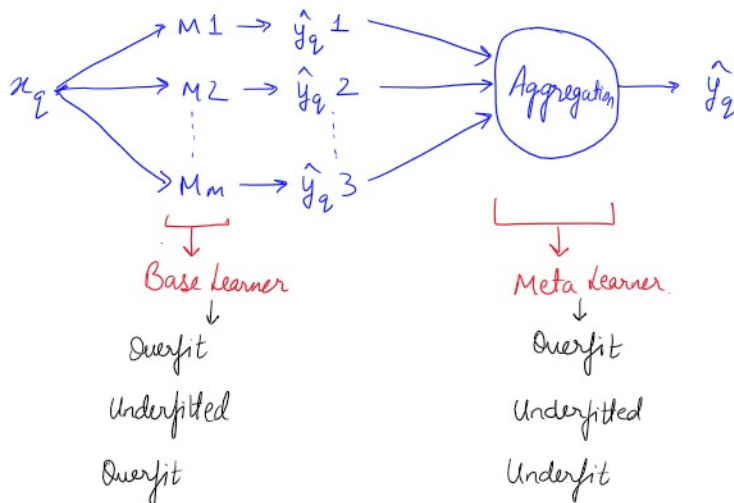
Prediction  $\rightarrow$



Base learner  $\rightarrow$  Same as Voting ensemble

Meta learner  $\rightarrow$  Classification  $\rightarrow$  Logistic Reg., D.T., N.B., etc.  
 $\rightarrow$  Regression  $\rightarrow$  Linear Reg., D.T., K-NN, etc.

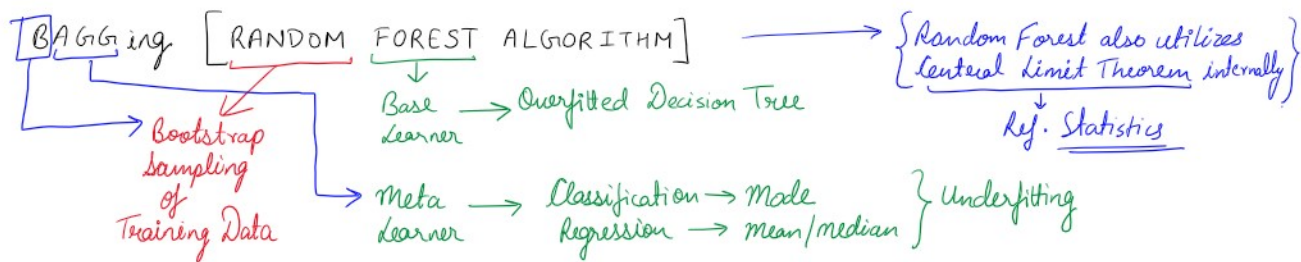
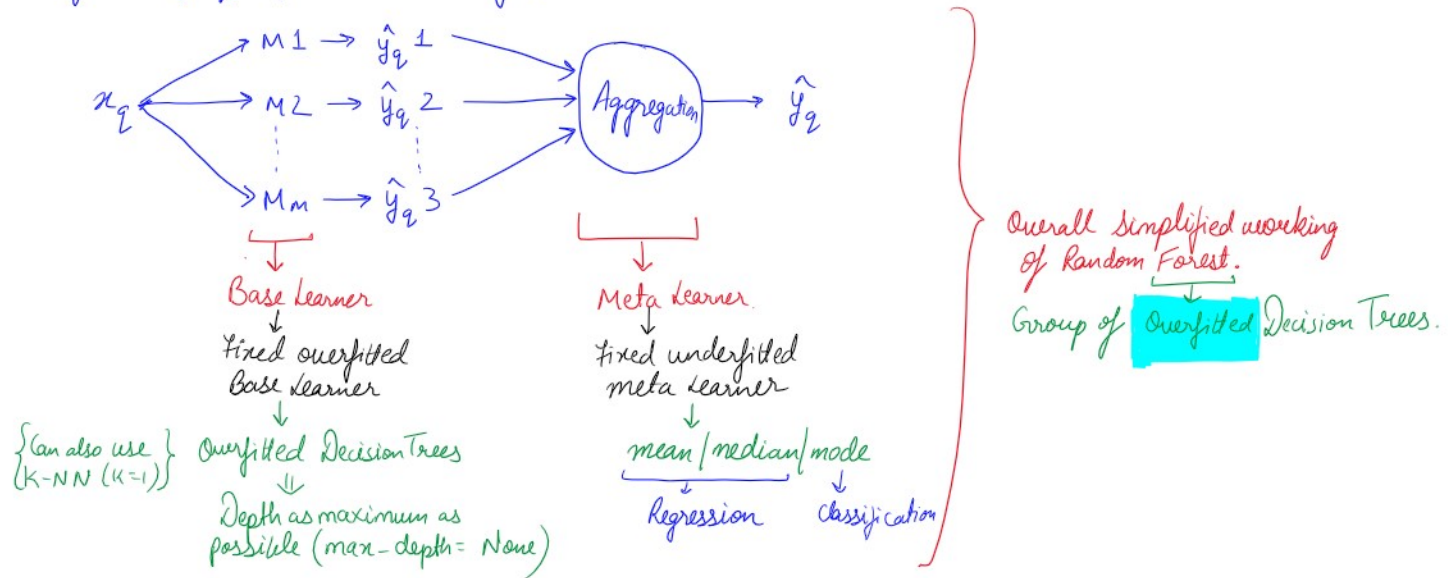
Now in Base learner as well as Meta learner, in both we can have overfitting or underfitting and as a result we have just ended up introducing more complexity than voting ensemble.



$\rightarrow$  Overfit X (Problem of Stacking)  
 $\rightarrow$  Underfit X (Stacking & Voting Ensemble)  
 $\rightarrow$  **Best fit** { Stacking  $\checkmark$   
Voting Ensemble  $\checkmark$  }  
But requires a lot of hyperparameter tuning

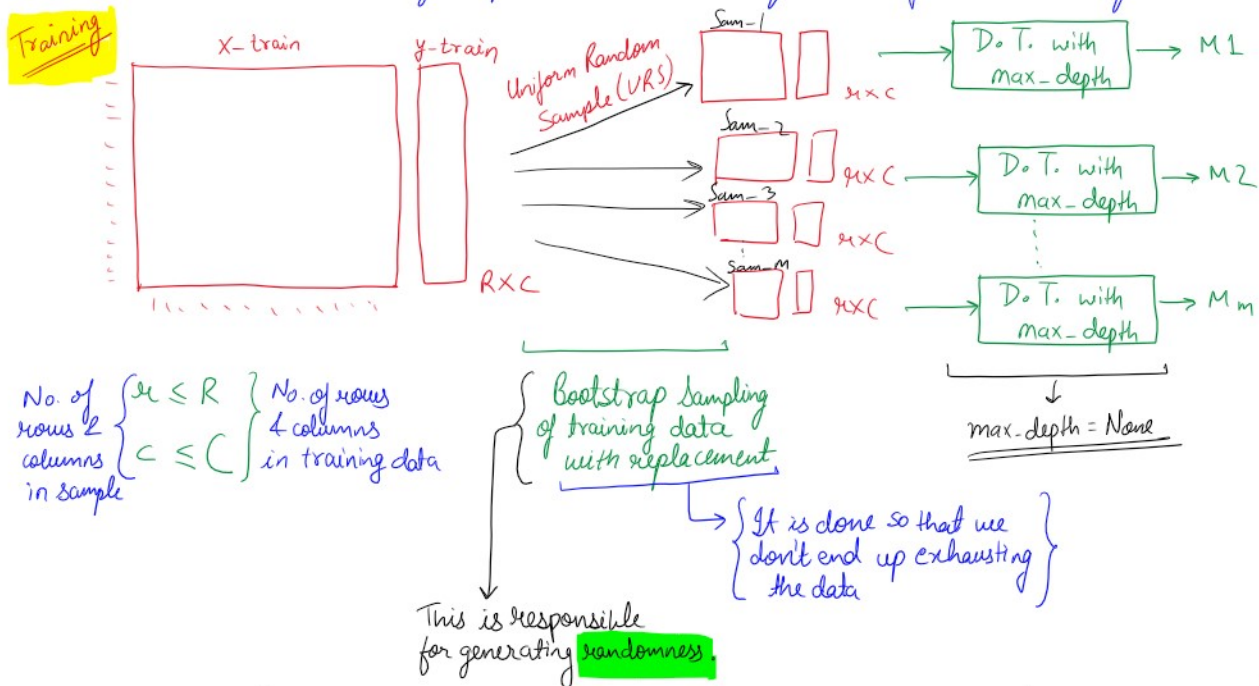
# # Bagging (Random Forest)

One of the best parts about bagging is that it is able to take care of the underfitting & overfitting problem without doing a lot of hyperparameter tuning (Little to None).



\* Bootstrap Sampling {The sampling is done on rows as well as columns}

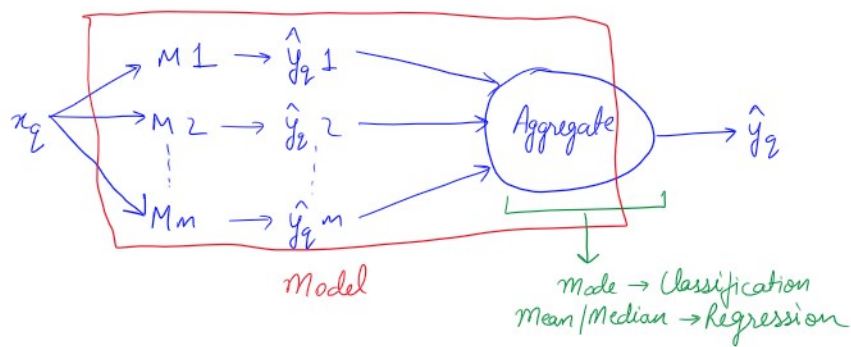
With Replacement → Any datapoint can be Randomly selected for Random No. of times.



\* All the samples will be different from one another as the Bootstrap Sampling will be done on rows as well as columns



Prediction



## Random Forest

- $\hookrightarrow$  Base learner  $\rightarrow$  Decision Tree (Overfitted)
- $\hookrightarrow$  Meta learner  $\rightarrow$  Mean/Median [Regression]  
mode [Classification]

Training Random Forest  $\rightarrow$  Bootstrap sampling with replacement for rows & columns of training data.  
 $\rightarrow$  On each sample  $\rightarrow$  Train an overfitted Decision Tree Model

Prediction with Random Forest

$\rightarrow$  Pass  $x_q$  to all models and get their predictions

$\rightarrow$  Aggregate all the predictions of model using: mean/median [Regression]  
mode [Classification]

[Code Example]