<u>Ensemble definition</u>:- A group of seperate things that contribute to a coordinated whole.

As the definition hints ensembling method use multiple week classifier to make a single powerful/robust classifier. Usually there are two types of ensambling method:-

1) <u>averaging method</u>:- make several weak estimator independently then average their prediction. It reduces the variance. e.g. Bagging, Random forest etc.

2) <u>Boosting method</u>:- make several weak classifier sequentially where one tries to reduce the bias of combined model. e.g. Adaboost, Gradient tree boosting etc.

## <u>Bagging</u>

Bagging stands for bootstrap aggregation.

<u>Bootstrap</u>:- Makes a random subset of dataset from the training dataset with replacement.

<u>Aggregation</u>:- average the result of all the weak classifier. Here every weak model will have the same priority during voting.

$$G(X) = \sum_{m=1}^{M} \frac{G_m(x)}{M}$$

Here, $X \rightarrow$ Data to be predicted.

$G_m(.) \rightarrow$ weak model.

$M \rightarrow$ Number of week model.

Out of bag estimation :- Each bootstrapped sample contains approximately 2/3 of the total training set. So we can use remaining 1/3 of traing data to calculate the error estimation, called out-of-bag error. If $M \rightarrow \infty$ then out of bag error gives an equivalent result to leave-one-out cross validation.

| Advantages | Disadvantages |
|---|---|
| Decrease variance | Increase bias |
| Better accuracy | Harder to interpret |
| Free validation set | still not additive |
| Support of missing value | More expensive. |

## Random Forest

It's a bagging technique which contains bunch of decision trees. Those are trained with the bootstrapped dataset. To make sure multipled tree

doesnot calculate the same thing we restrict the tree to choose the split between k feature out of n training feature. Value of k can be fine turned by using out of bag error. whichever k value will give us less out of bag error we will choose those.

We can make as many trees as we want using the same process. The trees should not have high deepth.

During testing we will average all the tree's result to get our final prediction..