

ENSEMBLE

METHODS / LEARNING

- * ensemble method use multiple base(0s) learning algorithms to give better performance

OR

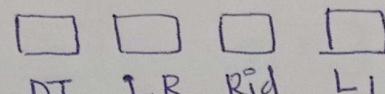
ensemble learning signifies the collection of multiple machine learning algorithms and combine them to form a good model is called as ensemble technique.

Imp:

- 1) All the basic models should be different, this can be done in two types

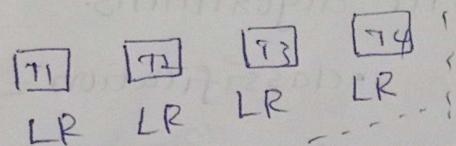


All the algorithms are different



↳ ensemble learning

giving different data for training every model



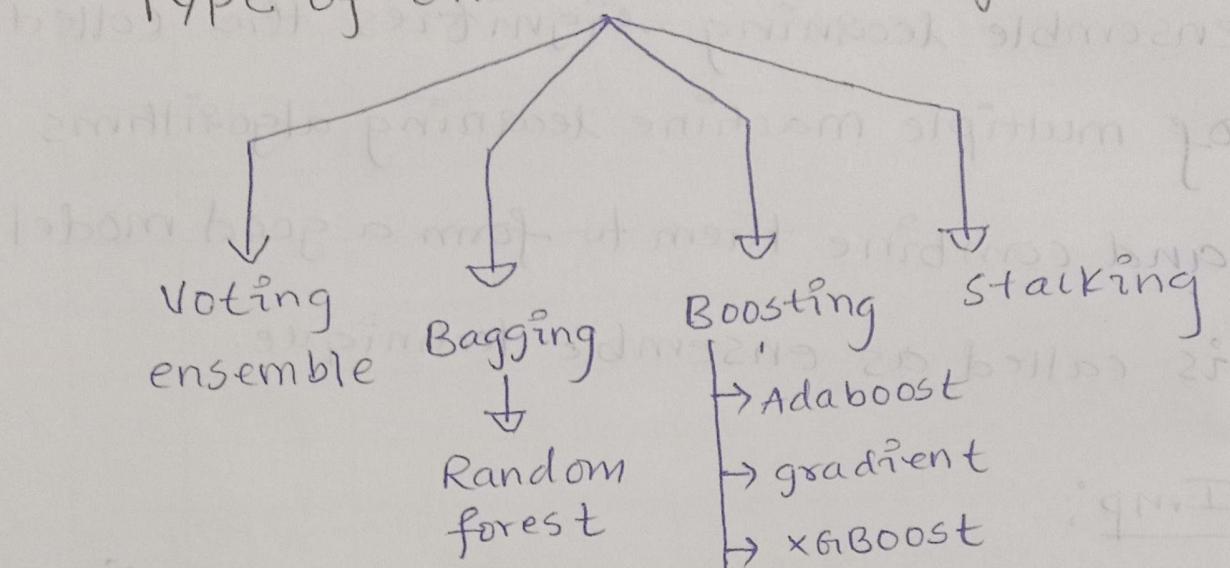
↳ ensemble method

71, 72, 73, 74 are different training datasets.

* For classification: The ensemble model gives the output based on majority count.

* For Regression: The out is based on the mean of all the outputs of every algorithm.

TYPE OF ENSEMBLE LEARNING



* The above types are existed because of two (or) 2 types of ensemble techniques (which are discussed on previous page).

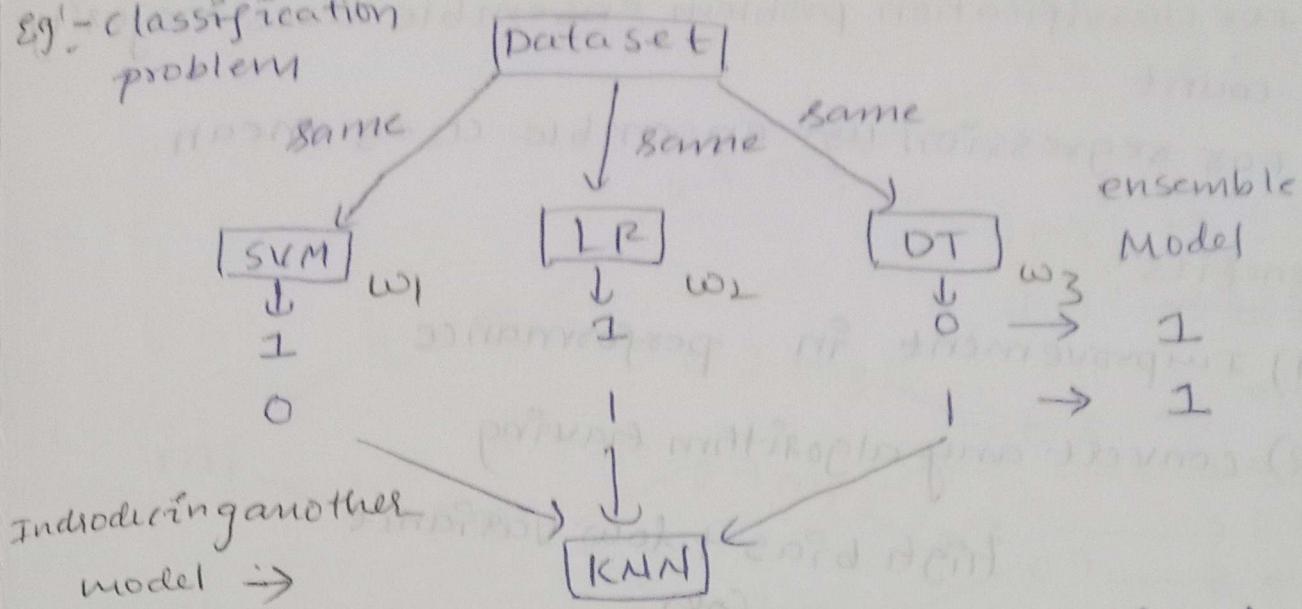
* Voting:- Based on ensemble model having different algorithms.

- classification = majority count

- Regression = Mean.

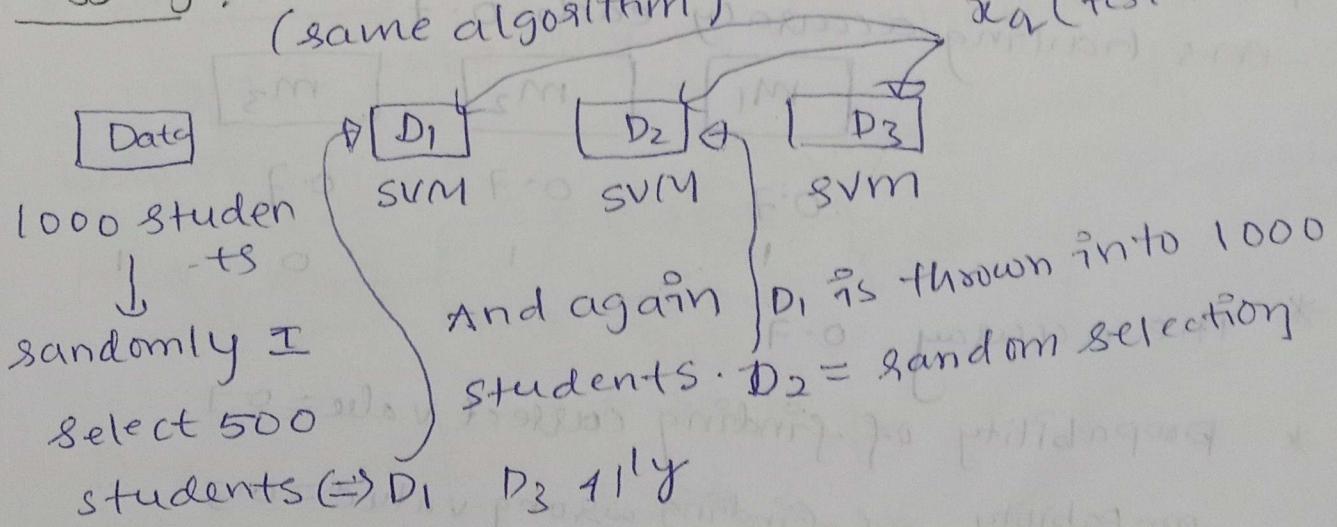
* Stacking:- Based on different algorithms

Eg:- classification problem



- * The KNN will train on the output or predicted values of different algorithms and assign weights based on that data. so that we can conclude either SVM, LR, or DT is performing well in the ensemble model.
- * Both voting and stacking are same but in stacking we additionally only about performance of each algorithm w.r.t ensemble model output by using another ml algorithm.

Bagging: Bootstrapped aggregation (Based on same algorithm)



- * For classification problem ensemble uses majority count
- * For regression use ensemble uses mean.

Benefits:

1) Improvement in performance

2) convert any algorithm having

High bias + low variance

(or)

low bias + high variance



Low bias + low variance

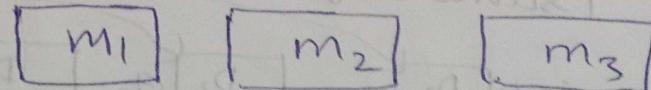
3) Robustness (?)

When to use :- Always

1. VOTING

1.1 \Rightarrow Intuition:-

problem :- If we have three different algorithms having accuracy each 0.7



0.7

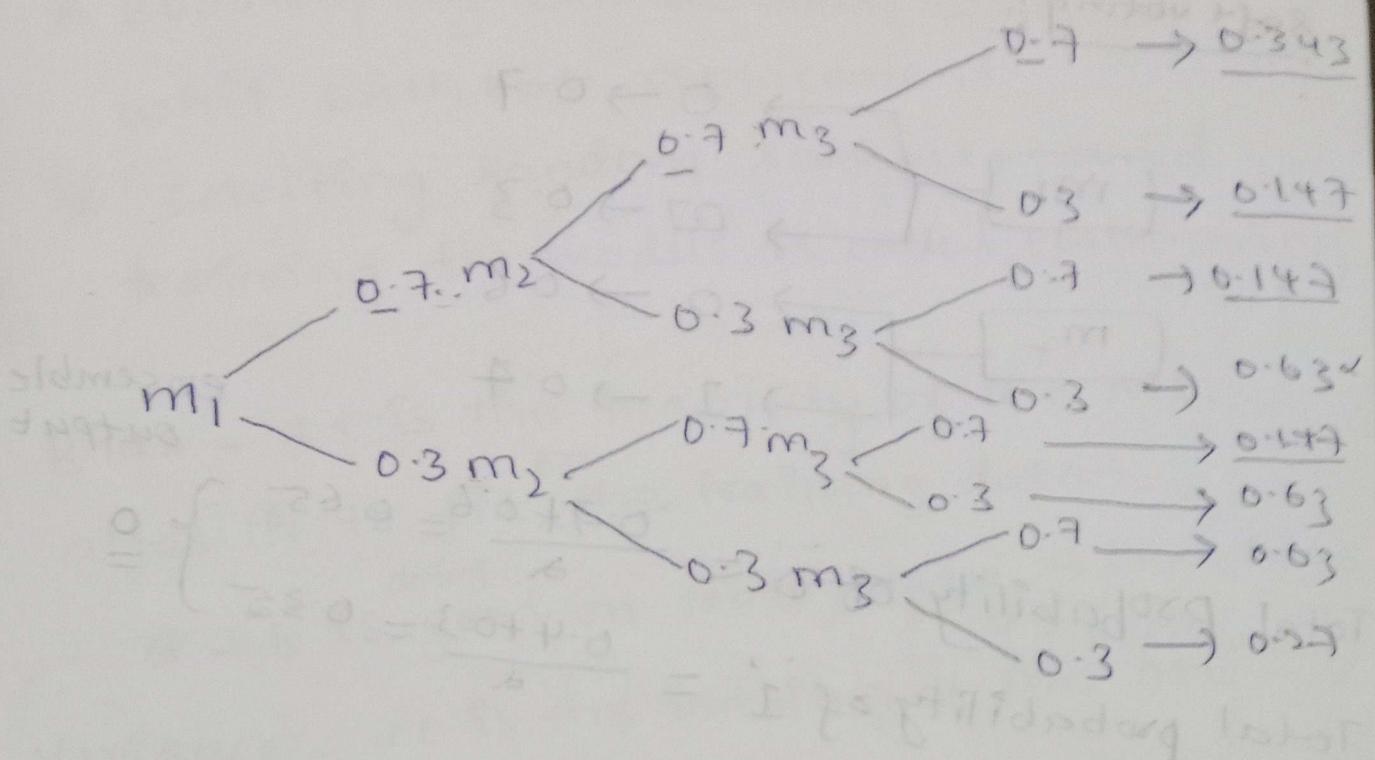
0.7

0.7

* The accuracy = 0.7

* Probability of finding correct value = 0.7

Probability of finding wrong value = 0.3



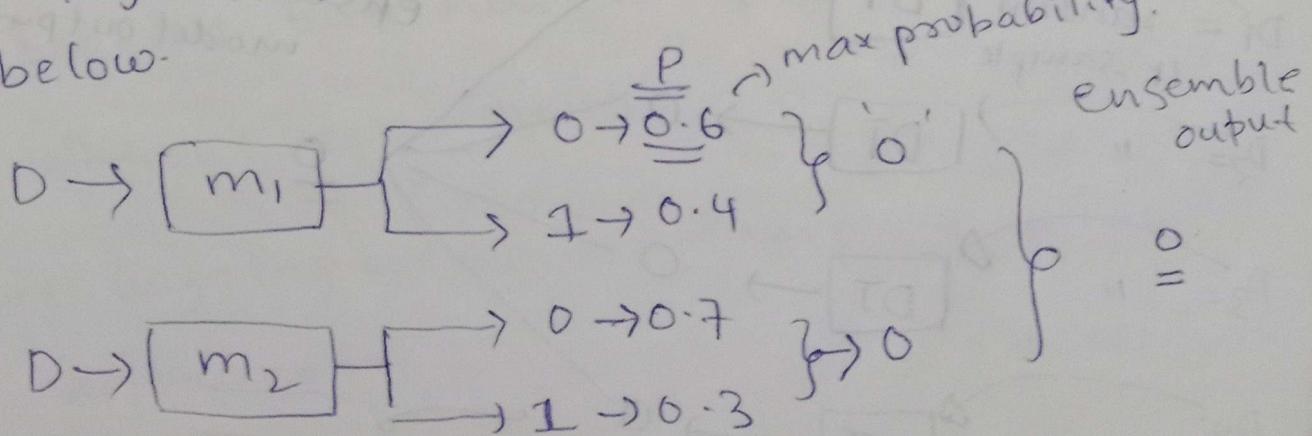
probability of at least two correct answers by three algorithms = 0.78

Voting classifier:-

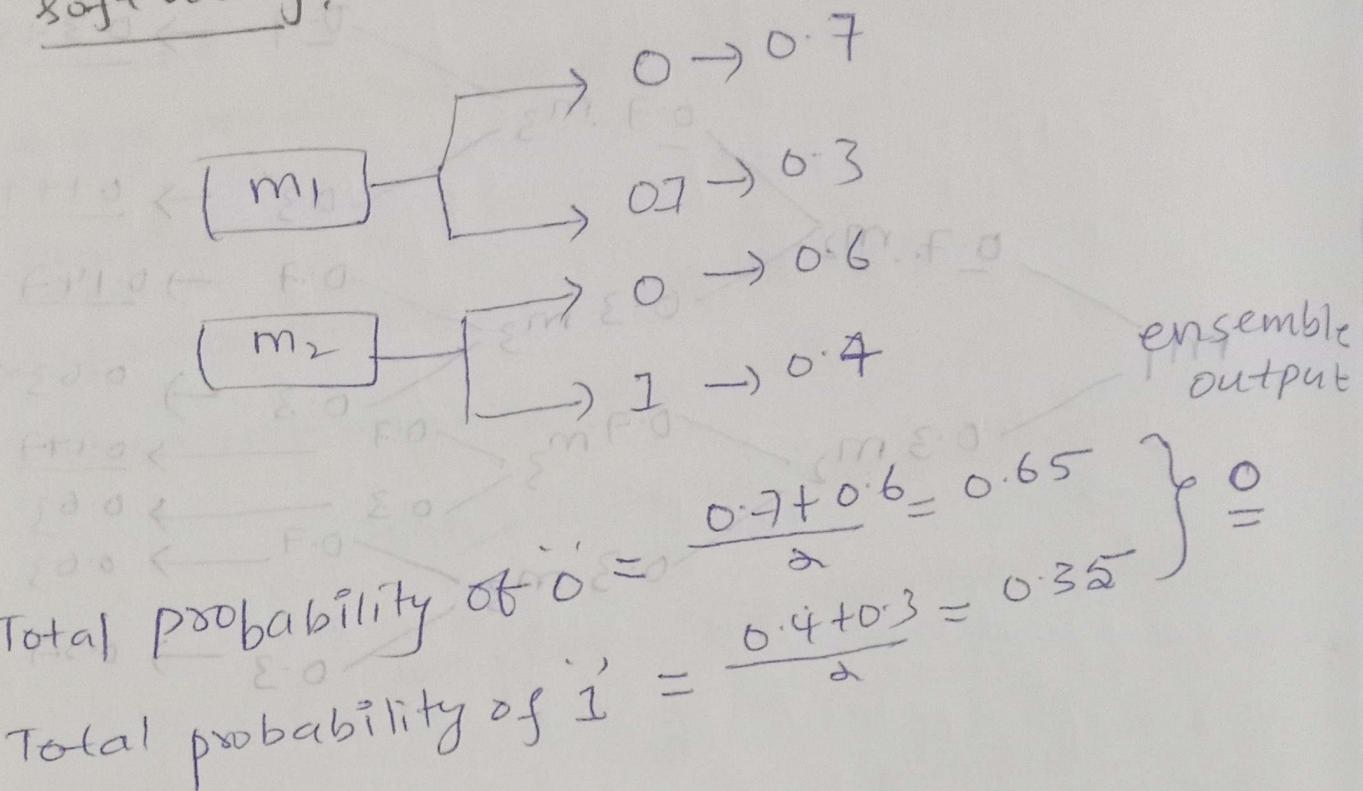
↓
soft voting

↓
hard voting

* hard voting is nothing but choosing the highest probability of the algorithm as shown below.



soft voting:

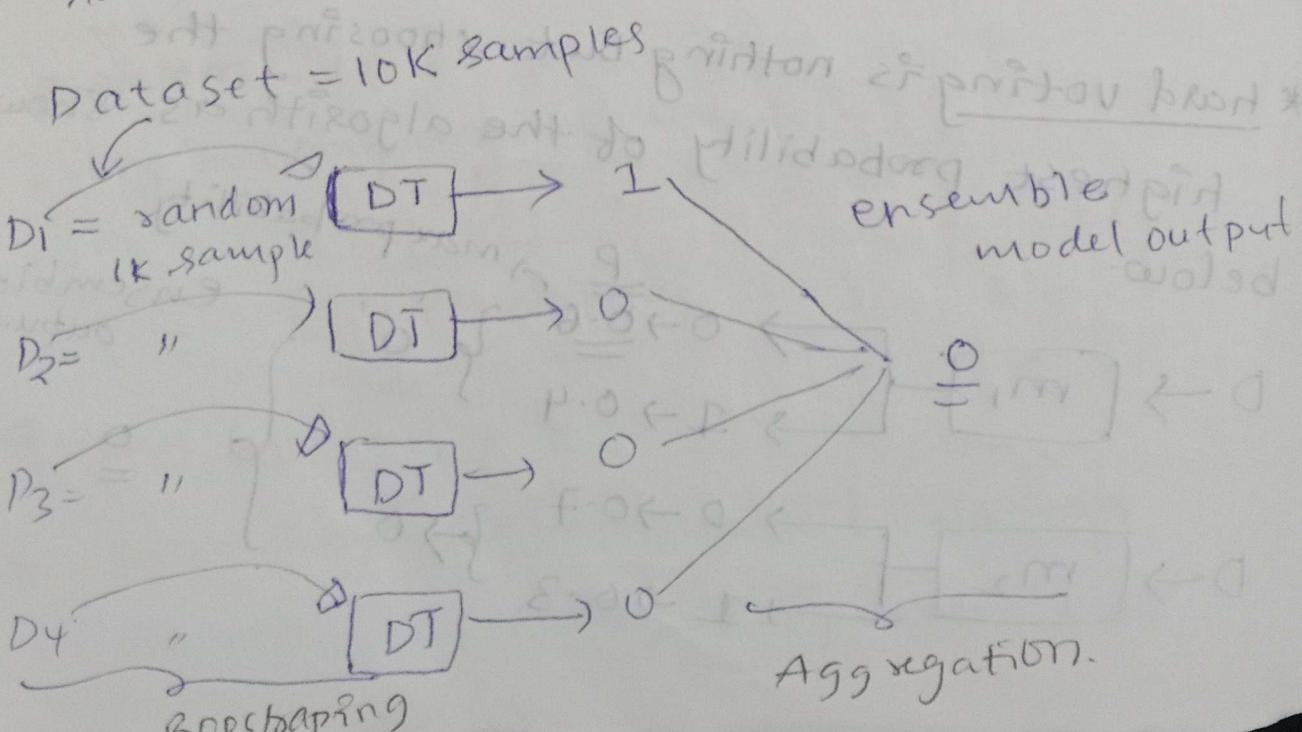


BAGGING :-

Bootstrapping

[from a given population selecting random samples]

Aggregation



- * Bootstrapping can be done with or without replacement.
 - * Selecting random samples and training the models is called bootstrapping.
- Bias:- It is the inability of the ML algorithm to fit on the training data doing good on training data \Rightarrow Low bias
- Variance:- It is the inability of the ML algorithm to fit on test data.
- * we use bagging so that the algorithm should maintain low bias and low variance. For ex: generally DT, SVM, KNN are HV type.
 - * Bagging gives consistent results.
 - * Pasting: Pasting is nothing but bagging but here we use bootstrap sampling without replacement.
 - * In Bagging we do sampling with replacement.

Random Subspaces: - Here we do column sampling with or without replacement.

Random patches: - Here we do row sampling and column sampling with or without replacement.

Random Forest :-

- Collection of trees
- In Bagging we do sampling (randomly)

Bagging Vs Random Forest

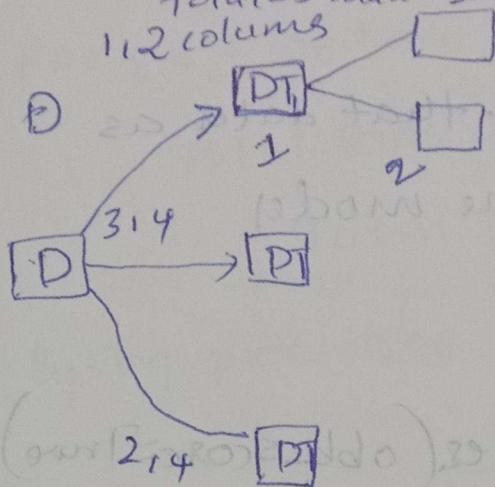
1. In bagging we can select any base estimator like DT, logistic regression, KNN, SVM - or any. But in random forest we by default the estimator is decision tree.
2. This difference is based on feature sampling.
 - In bagging only tree level column sampling will happens.
 - & In bagging & RT we are using same DT. If it's bagsins has no difference with RF.

Bagging

feature selection = 2

Total = 5 column

1, 2 columns



Here the decision

tree forms

nodes only on 1, 2 column

when we consider DT1

(not like each decision

tree we can observe
other columns also, however
we given only 1, 2

columns

Random forest do feature selection is

Random forest do feature selection is
node level column sampling.

*) Due to this randomness of the random forest

it will be more

OOB SCORE: out of bag

0.8

Evaluation: bootstrap = True (with replacement of rows).

* Generally when we are using Bagging or Random forest algorithm, with replacement

mathematically proven that 37% of the data
is unseen by the algorithms.

Also by using oob we can use that data as
a validation set to test the model.

Syntax:

```
rf = RandomForestClassifier(oob_score=True)
```

```
rf.fit(X_train, y_train)
```

```
rf.oob_score_
```

output = 0.785 (For example)

Feature Importance: Finding the importance of

the column w.r.t Random forest or any algorithm.

* In this way we can reduce the columns and
also feature importance gives interpretability.

For example: In a classification problem whether
loan will get sanctioned for not.

* For someone whose loan is not sanctioned
he reaches out to me and asks why?

* Then as I know the important columns I can
interpret them.

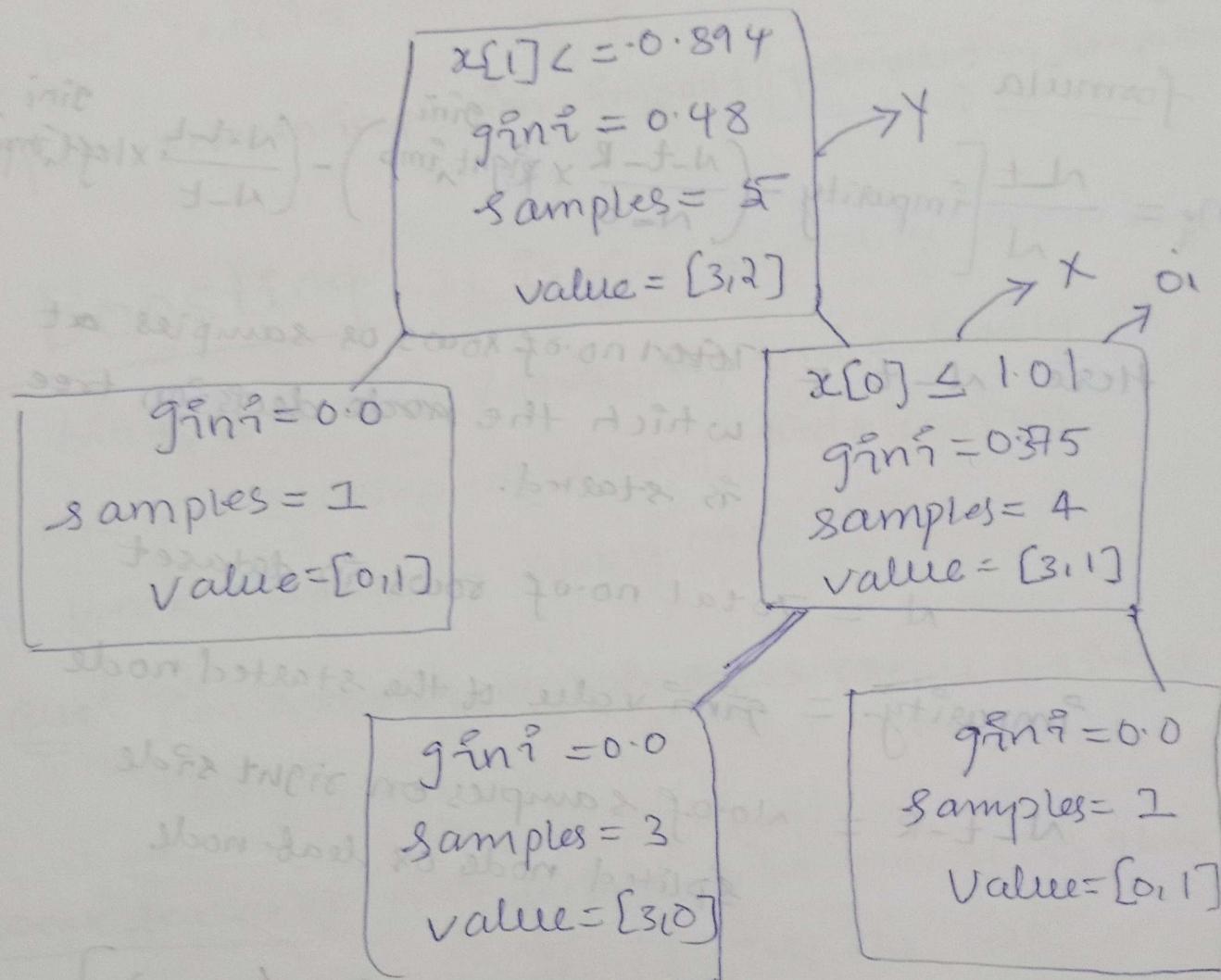
How feature is calculated :-

for ex:- we have

Total rows or samples = 5

Total columns = 2

Using plot-tree, we get after fitting with total rows



column importance of 0th column

No. of nodes ^{that} are splitted on 0th

$$= \frac{\text{total no. of nodes}}{\text{total no. of nodes}}$$

= $\frac{\text{Importance of nodes that are splitted on 0}^{\text{th}}}{\text{Total Importance of nodes}}$

$$\text{For } 0^{\text{th}} = \frac{x_i + y_i}{x_i + y_i} \left\{ \begin{array}{l} \text{for our example} \\ x_1=1, y=1 \\ \therefore \text{total sum} = 1+3=4 \end{array} \right.$$

$$1^{\text{st}} = \frac{y_i}{x_i + y_i} \quad \begin{array}{l} \text{sums to same total} \\ \text{sums to same total} \end{array}$$

Now how feature importance of nodes are calculated.

formula:

$$n_i = \frac{n_t}{N} \left[\text{impurity} - \left(\frac{n_t - R}{n_t} \times \text{right imp} \right) - \left(\frac{n_t - L}{n_t} \times \text{left imp} \right) \right]$$

Here n_t =

$$280 = 280$$

$$4 = \text{samples}$$

$$[1,8] = \text{sub set}$$

N = Total no. of rows or samples at which the node decision tree is started.

$$I = 280$$

$[1,8] = \text{subset}$

N = Total no. of rows in dataset

n_t = Total no. of rows in node

$R = 20$ = No. of samples on right side

$L = 10$ = No. of samples on left side

$n_t = 280$ = Total no. of samples

$n_t = 280$ = Total no. of samples

$n_t = 280$ = Total no. of samples

$$* (n)_{0^{\text{th}}} = \frac{4}{280} \left[0.375 - \left(\frac{1}{4} \times 0 \right) - \left(\frac{3}{4} \times 0 \right) \right]$$

$$x = (n)_{0^{\text{th}}} = \frac{4}{5} \times 0.375 = 0.3$$

$$* y = (n)_{1^{\text{st}}} = \frac{5}{280} \left[0.48 - \left(\frac{4}{5} \times 0.375 \right) - \left(\frac{1}{5} \times 0 \right) \right]$$

$$= (0.48 - 0.3) = 0.18$$

$$\text{O/F} = \frac{x}{x+y} = \frac{0.3}{0.3+0.18} = \frac{0.3}{0.48} = 0.625$$

$$1 - \text{O/F} = \frac{0.18}{0.48} = 0.375$$

* If the data has high cardinality then we should then use `PermutationImportance` from `sklearn.inspection` import

Boosting: It is for primarily reducing bias, and also Variance, and a family of machine learning algorithms that convert weak learners to strong ones.

Ques: Can a set of weak learners create a single strong one?

weak learner: It is defined to be classifier that is only slightly correlated with the true classification. means an algorithm which has accuracy just more than 50%.

(or)

refers to simple models that do only slightly better than random guessing.

Adaboost:

1) weak learners

2) decision stumps

3) -1 and +1 (instead of 0 & 1)

Decision stumps: Actually these are types of weak learners. It is a type of decision

$\boxed{\text{max-depth} = 1}$

Decision tree in which it's [means split happens at one node only]

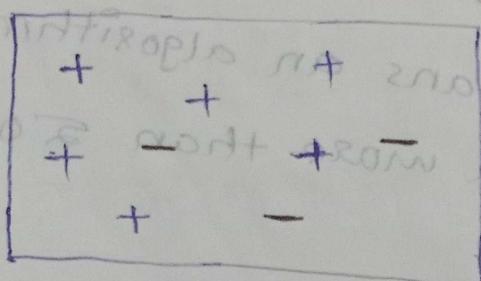
* We use decision stumps as our weak learners in the Adaboost because we can see what's going on in the algorithm and

for good results

* Adaboost is a stage wise additive method

(It is formed by multiple weak learners). Here

we add weak learners sequentially



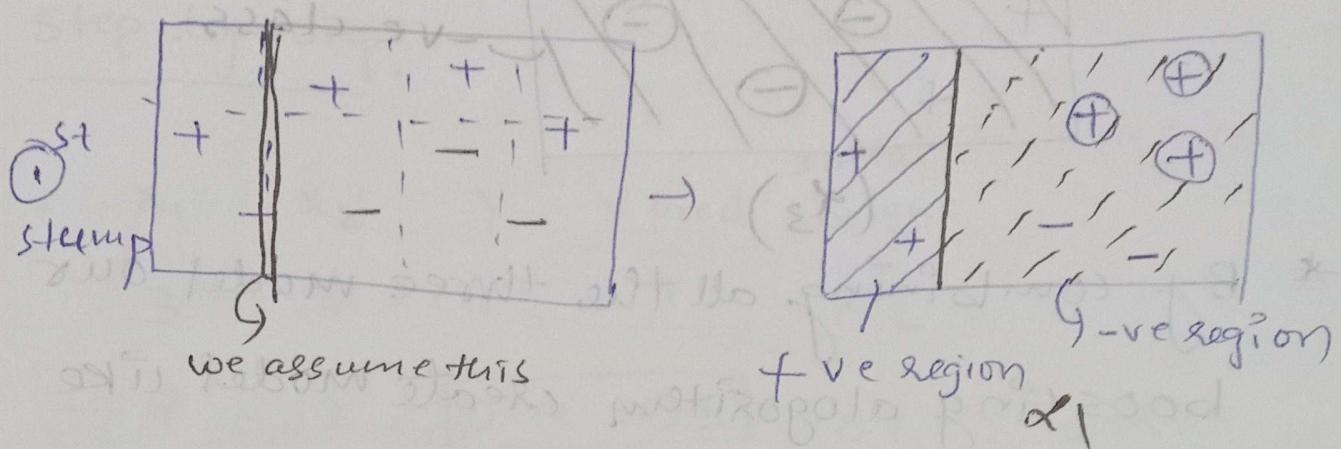
+ → got placement

- → not got placement

- This is our data set is based on CGPA,

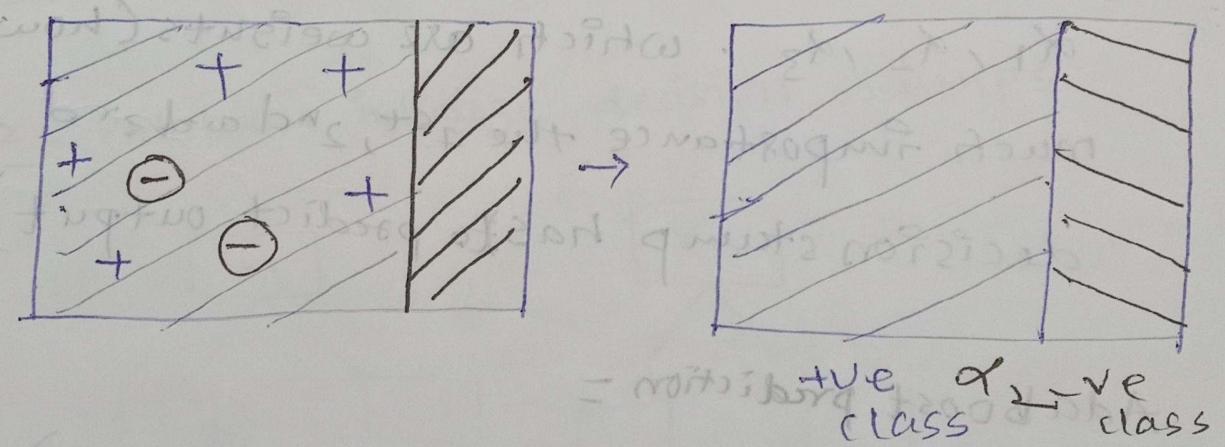
IQ and output = placement + not placement

* Now we apply random decision stumps and select one model in such a way that entropy is minimum.



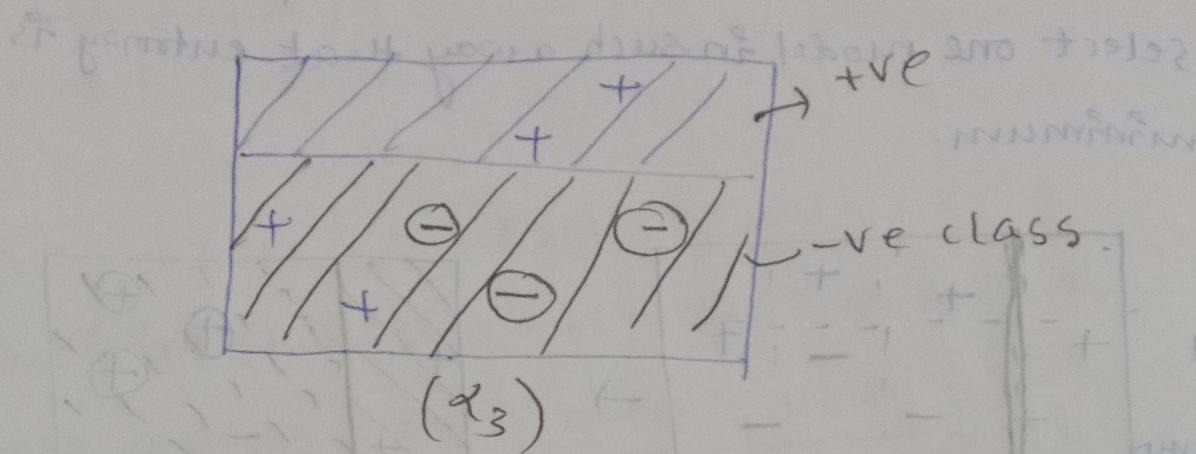
② By analyzing above results three \oplus points are incorrectly classified.

Aim!: now our 2nd stump aim is create a model in such a way that 3 \oplus points are correctly classified. This Aim is also called as **upsampling**

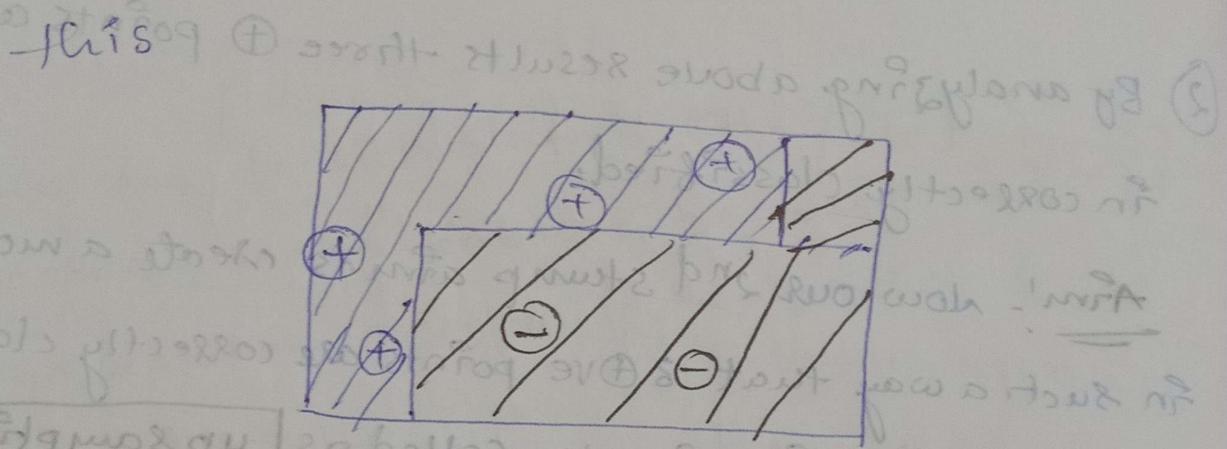


3rd Stump - Aim! It should create a model in such a way that the 2 \ominus points should correctly classified.

base squate, mohit probhab piggas sarewala *



* By combining all the three model our boosting algorithm create model like



* In the same process we also calculate $\alpha_1, \alpha_2, \alpha_3$. which are weights (how much importance the 1st, 2nd and 3rd decision stamp has to predict output).

Adaboost prediction =

$$\text{sign}(\alpha_1 \times h_1(x) + \alpha_2 \times h_2(x) + \alpha_3 \times h_3(x))$$

$$\left\{ \begin{array}{l} \text{if } \sum_{i=1}^n \alpha_i h_i(x) > 0 \text{ then } \Rightarrow +1 \\ \text{if } \sum_{i=1}^n \alpha_i h_i(x) < 0 \text{ then } \Rightarrow -1 \end{array} \right\}$$

here $h_1(x), h_2(x), \dots$ are the decision stumps.

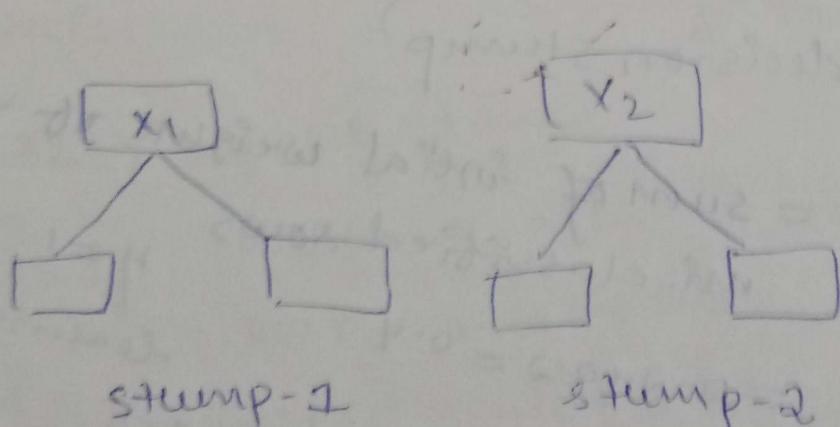
Step-by-step:

x_1	x_2	y	y_{pred}	Initial weight
3	7	1	1	0.2
2	9	0	0.1	0.2 } mis
1	4	1	0	0.2 } classified
9	8	0	0	0.2
3	7	0	0	0.2

Step-1: our adaboost model initialize weights to each and every row by using formula = $\frac{1}{\text{no.of rows}}$

$$\text{In above case } w = \frac{1}{5} = 0.2$$

Step-2: Adaboost initialize decision stumps



* In the above two stumps which has less entropy or gini that decision stump (weak learner is selected).

Step-3 Now we predict (x_1, x_2, \dots, x_n) in \mathcal{X}

The main aim of adaboost is to assign weights to each and every decision stump.

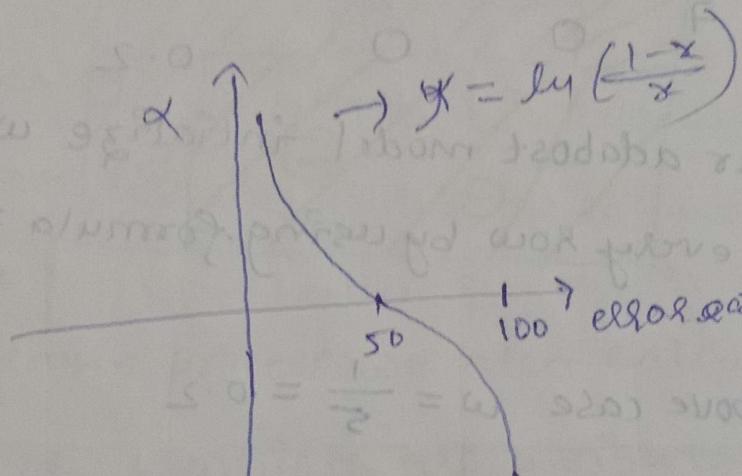
* By using one decision stump we calculate the y_{pred} values

$$\eta_d = \frac{1}{2} \ln \left(\frac{1 - \text{error}}{\text{error}} \right)$$

learning rate

Ex: If error rate = 0%, the η_d is max
 $\text{error rate} = 100\%$, the η_d is min

$$\eta_d = 50\% \quad d = 10$$



m_1 (1st decision stump)

error = sum of initial weights of the misclassified rows

$$= 0.2 + 0.2 = 0.4$$

$\eta = 1$
learning rate

$$\eta_{d_1} = \frac{1}{2} \ln \left(\frac{0.6}{0.4} \right) = \frac{1}{2} \ln \left(\frac{3}{2} \right)$$

$$= \frac{1}{2} (0.477 - 0.301) = \frac{0.176}{2}$$

$$w_2 = 0.69$$

$$w_3 = 1.09$$

$$\alpha_1 = \frac{0.405}{a} = [0.20 = d_1]$$

Aim:- To upsample misclassified weights.

For misclassified:-

$$\text{new-wt} = \text{curr-wt} \times e^{\alpha_1} = 0.2 \times e^{0.2} = 0.24$$

For correctly classified:-

$$\text{new-wt} = \text{curr-wt} \times e^{-\alpha_1} = 0.2 \times e^{-0.2} = 0.16$$

* Aim) To get the sum of updated weights is 1

as initial weights we

normalize them by dividing

each one with Total sum

Now our data set looks

like this:

updated	normalized
0.16	0.166
0.24	0.25
0.24	0.25
0.16	0.166
0.16	0.166
<u>0.96</u>	<u>1</u>

$x_1 \ x_2 \ y \ \text{newwt} \ \text{range}$

0.166 0 - 0.166

0.25 0.166 - 0.416

0.25 0.416 - 0.666

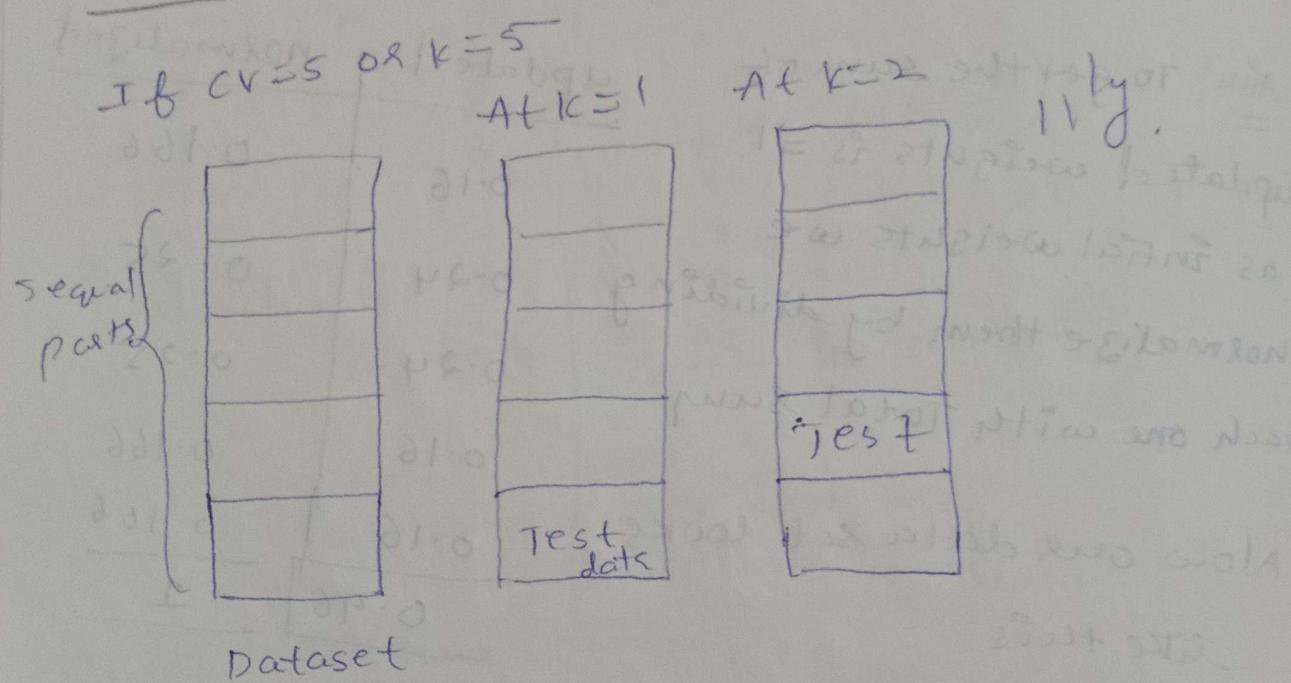
0.166 0.666 - 0.832

0.166 0.832 - 1.0

* Now add bias & select 5 random numbers from 0 to 1
For eg. these are those numbers :-

0.13 → 1st row
0.43 → 3rd row
0.62 → 3rd
0.50 → 3rd
0.8 → 24th

Cross-validation :- or K-fold CV



learning rate hyperparameter (η) :-

If η is small then there is a chance of underfitting because the upsampling weights are not have much difference with weights of correctly classified.

Bagging Vs Boosting:

1. Type of model :-

→ In bagging we use a model which

has "Low Bias & High Variance"

e.g.: fully grown decision tree.

→ In boosting we use model which has

"High Variance & Low bias"

e.g.: decision stump.

2. Sequential learning vs parallel :-

→ bagging → models will train parallelly at a time

→ boosting → model will train one after the other in sequential order.

3. weightage of base learners:-

→ Bagging = In bagging all the base model will get or have equal priority means equal weights.

→ Boosting = In boosting every model will get assigned by different weights