

Ensembles Learning

- 1) In ensembles learning we look at multiple classifier & combining the output of the multiple classifier in order to get better prediction or classification accuracy.
- 2) In certain conditions, where the classifier outputs are independent of each other and make error in an independent manner it is possible that by combining the outputs of the multiple classifier in order to get better prediction or classification accuracy.

Ensembles classifier

- 1) In this we have multiple learner or classifier which give multiple output which may same or all different, or some are same & some are different.
- 2) Multiple learner will give different decisions if we want to combine a decision output by the multiple learner.

For this we have to generate a group

of base learner.

- So How to generate a group of learners.
- 1) use Different Algorithm
 - 2) Different learner use different hyper parameters.
 - 3) Or Same algorithm but different ^{values of} parameter you get different models as output will be different. They may have ⁴⁾ different representation
 - 5) different training sets

For ex:- 1) In ANN we can select different no of layer, weight, bias, neuron as hyper parameter
2) In decision tree the strategy for deciding this attributes like entropy, gini index. to modify the decision tree algorithm.

→ For ex 2 class - classification

Error of each learner is E
if $E \geq 0.5$

means may be accuracy is not high but error should be $E \leq 0.5$.

* Different base learner satisfy certain independence conditions, then it is possible to combine weak learners, each having error ~~greater~~^{less} than 0.5 in order to get a very strong learner.

* Importance of Bias and Variance.

Bias:- Is a measure of how flexible the model is. So if the model is very flexible or very powerful then the bias is low.

Variance:- Variance is high when if you give different subsets of data as training set, the models output are very different then we say variance is high.

* So in machine learning our objectives is reduce both bias & variance.

* So using ensembling learning we want low bias & variance.

* Ensembles use multiple trained model

1) If we assume this models have low bias & high variance.

By combining the output, we can get low variance while maintaining low bias.

Even at initial individual hypothesis may have high bias by combining models we get classifier, which has low bias.

2) So by ensembles if we get low bias & low variance we will get less overfitting no need worry about stopping criteria.

* Combining weak learners

→ But error cannot be higher than 0.5;
error must be less than 0.5.

$$E < 0.5$$

So to combine weak learner following are conditions.

- 1) Assume n independent learner.
- 2) Each has accuracy $\neq 0\%$ or each with error 0.3 or 0.7 accuracy.

3) If all n learners has same output
 If all the learners agree on the output
 you are getting higher confidence.

For ex. - Ten learners with accuracy 0.7

$$\therefore \text{Confidence} = 1 - (1 - 0.7)^{10} = 1 - 0.3^{10} \text{ close to } 1.$$

* How to Combine the learner.

→ Combination can be done by voting.
 based on majority combination or majority voting

→ Voting can be given based on weight.

1) Unweighted Voting

2) Weighted Voting.

Weight \propto accuracy

Weight \propto $1/\text{variance}$

$$y = \sum_{j=1}^M w_j d_j \quad d_j \text{ is o/p of classifier} \\ w_j \text{ is weight of classifier}$$

$$\text{Each } w_j \geq 0, \quad \sum w_j = 1$$

ENSEMBLE TECHNIQUES

Two types of Ensemble techniques

1) Bagging - (Bootstrap Aggregation)

1) Random Forest

2) Boosting

1) ADABoost

2) Gradient Boosting

3) XgBoost

multiple learner

1) Bagging Ensemble

Base model

data

D^m

D^{m_1}

D^{m_2}

D^{m_3}

D^{m_k}

m_1

m_2

m_3

m_k

$m \leq n$

sampling

final output based on majority voting

Row Sampling with Replacement

* Ensemble means combining multiple models

done by combining multiple models

* In bagging from a dataset having n samples, select D^m dataset & provide to each model m_1, m_2, \dots, m_k with row sampling with replacement

* Resample data & give it other model

- * So each model is trained with different data set.
- * Then using test dataset each model is tested for binary classifier & based on voting of majority final OIP is selected.

* In Bagging another name is Bootstrap because it selects multiple datapath for training multiple machine learning model & finally Aggregation for voting based OIP classifier.

* For example in Random Forest tree we use multiple decision tree.

In single decision tree when we create trees with high depth we get Low Bias but High Variance

* But in Bagging Approach we use data with sampling with replacement we use data sampling & feature engineering & apply to multiple learner machine model & we take aggregation output based on

voting due to which we get low variance.

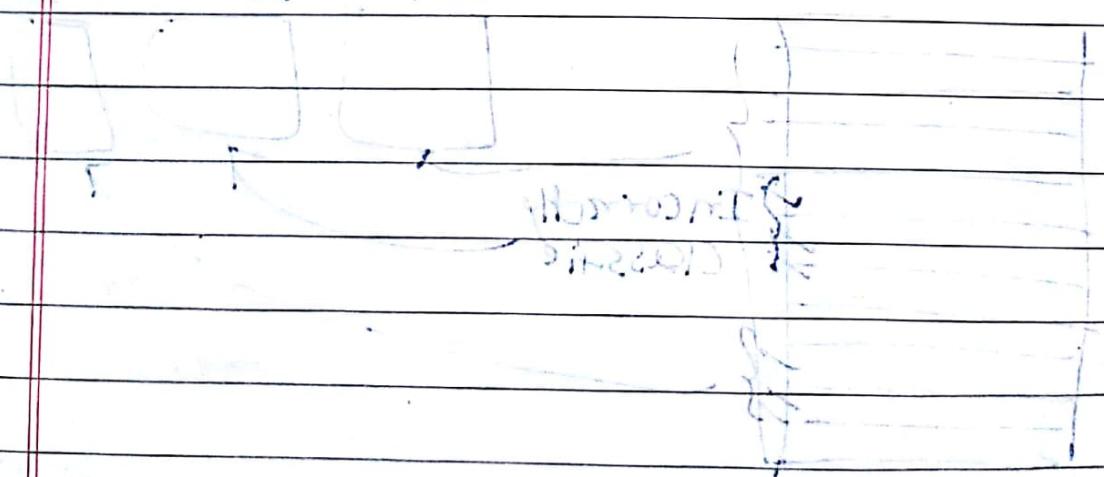
This multiple model become expert on different data set.

while getting high variance to low variance we are not depend on only one model output we take majority voting so get low variance.

* So in classifier uses majority vote

* In regression it uses mean or median as output.

* Hyperparameter : No of decision tree in random forest.

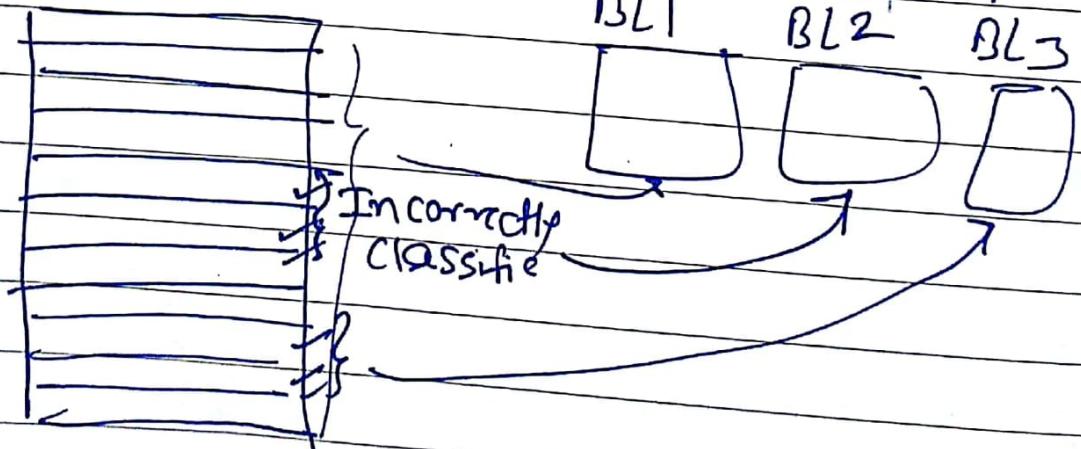


Topic on Boosting Techniques: (ADBoost)

* Boosting

- 1) Consider Data set: D
- 2) Base learner: BL1, BL2, BL3 ... BLn & selected
- 3) From the dataset give all example to first base model & check the accuracy.
- 4) Those example from dataset gives missclassified, Again give that example to another base model BL2
- 5) And for BL2 those example are missclassified give to BL3 & so on
- 6) It will go on unless & until we created base model with good accuracy.
- 7) Base model are created sequentially.

* ADABOOST



* ADA BOOST

- 1) Consider feature f_1, f_2, f_3 OIP
- 2) with multiple feature training example
- 3) Consider all training example have some sample weight

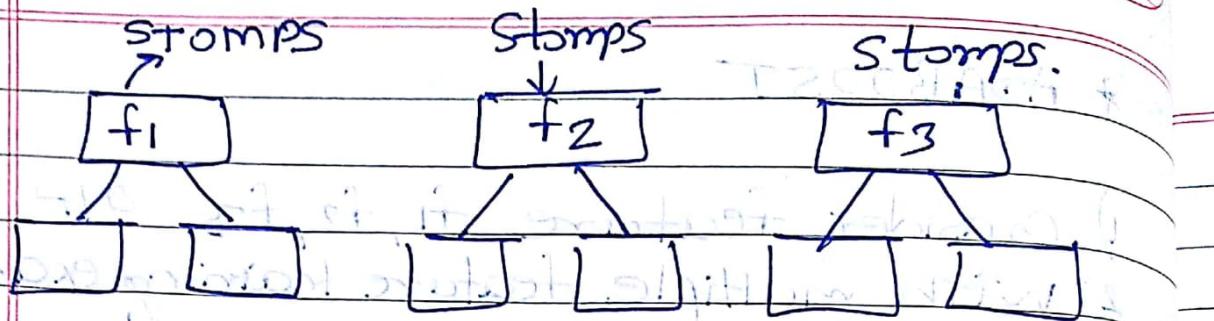
Initial Samp weight $w = 1/n$ n is no. of training example. (no. of records)

A) Initially all record assign same weight for ex.

- Step 1
- 1. f_1, f_2, f_3 OIP, sample weight $n=7$
 - 2. Assign weight $w = \frac{1}{7}$
 - 3. Assign weight $\frac{1}{7}$ to all record.
 - 4. No $\frac{1}{7}$ to all record.
 - 5. Yes $\frac{1}{7}$ to all record.
 - 6. Assign weight $\frac{1}{7}$ to all record.
 - 7. Assign weight $\frac{1}{7}$ to all record.

Step 2), Create a first base model BL1 with decision tree if + abr.

- 2) Create a base learner sequentially,
- 3) Create a decision tree only with one depth
- 4) Consider feature f_1 & create decision tree with one depth called STOMPS



5) From this three Stomps we have to select one decision tree as base learner model.

6) How to select first base learner from above three Stomps
→ Using Entropy or Gini Index.

Those decision tree entropy is less
Select that decision tree as base learning model out of above 3 Stomps one with less entropy.

Step 5) Once we select one decision tree

1) from above Stomps which have less entropy, check how many records it gives correct output

for ex → f_1 ← with less Entropy

Correctly classified

i-Incorrectly classified more balanced

2) Find out Total error = $\frac{\text{No. of incorrectly classified OIP}}{\text{Total No. of OIP}}$

for example in this table.

$$\text{Total error} = \frac{7}{10}$$

i :- is incorrectly classified OIP.

7 :- total no. of output not +

Step 4) Find out Performance of Stump.

$$\text{Performance of Stump} = \frac{1}{2} \log_e \left(\frac{1 - TE}{TE} \right)$$

Performance of Stump means how the stump is basically classified.

$$\text{Performance of Stump} = \frac{1}{2} \log_e \left[\frac{1 - \frac{1}{7}}{\frac{1}{7}} \right] = 0.896$$

→ Total error of performance of Stump is calculated to update the weight of data record.

So in boosting we need to pass wrong classified record to next Stump.

So for the same we need to update weight in such a way that for wrong classified record increase the weight & for correctly classified record decrease the weight.

Step 3) Updation of weight

Ex: use following.

$$\text{Performance of stump} = 0.895 \quad T.E = \frac{1}{7}$$

- 1) Update the weight for incorrectly classified record OIP.

$$\text{New Sample weight} = \text{Old weight} \times e^{\text{Performance of Stump}}$$

$$= \frac{1}{7} \times e^{0.895} = 0.349$$

so weight increases

- 2) Update the weight for correctly classified record OIP

$$\text{New Sample weight} = \text{Old weight} \times e^{-\text{Performance of Stump}}$$

$$= \frac{1}{7} \times e^{-0.895} = 0.05$$

so weight reduces

So using weight updation we update weight

- for incorrectly classification weight is increase.
- for correctly classification record weight is decrease.

for ex - Updated weight is, ^{update} weight

f_1	f_2	f_3	O/P	Old weight	New weight	Normalised weight
brown	brown	brown	Yes	1/7	0.05	0.07
brown	brown	brown	No	1/7	0.05	0.07

Initially we take Yes $\rightarrow 1/7 = 0.05 + 0.07$

Yes $\rightarrow 1/7 = 0.05 + 0.07$

No \leftarrow incorrect classified $\rightarrow 1/7 = 0.34 + 0.513$

Initially brown yes $\rightarrow 1/7 = 0.05 + 0.07$

After Normalized $\rightarrow 1/7 = 0.05 + 0.07$

Initially green green $\rightarrow \sum = 1 - \sum = 0.68$ ($\sum = 1$)

As we take summation of old weight = 1

But summation of New weight is not equal to 1

∴ Normalised the weight for that.

$$\text{New weight} = \frac{\text{of New weight}}{\sum \text{ of New weight}} = \frac{0.05}{0.68} = 0.07$$

Use this Normalized weight.

Step 6:- Create New database Bucket.

by the range.

$f_1 f_2 f_3 \cup P$ Normalised Bucket. So by dividing

Yes $0.07 \rightarrow 0.007$

No $0.07 \rightarrow 0.007 + 0.14$

Yes $0.07 \rightarrow 0.007$

Yes $0.07 \rightarrow 0.14 + 0.21$

wrong
classified
record.

$\rightarrow 0.0513$

Yes 0.07

No 0.07

the weight into
range we

create a
bucket

& select the
records in

of which we select wrong
classified record.

i) Create different iteration on selected random value

ii) for ex 0.43 which is in the range of bucket 0.21 to 0.727 select the record which falls in this particular range.

iii) In this range wrong classified record is also present.

iv) Select this record, take new value of iteration 0.31 select the bucket & select new record values which are in this bucket

v) In this way create a new dataset

$f_1 f_2 f_3 \cup P$

- 1)
- 2)
- 3)
- 4)

\therefore take a random no. \Rightarrow it's in

vi) Use this new dataset for creating new decision tree Stumps.

Step 7) Note: Repeat the process again from Step 2, i) Select decision tree based on less entropy value.

ii) Calculate correct & incorrect classified record

iii) find Total error & performance of Stump.

iv) Update weight.

v) Create bucket based on range.

vi) By iteration select record from bucket & create new feature table

& use new decision stumps, find error again.

Repeating this process again & again we get very less error which we find at start.

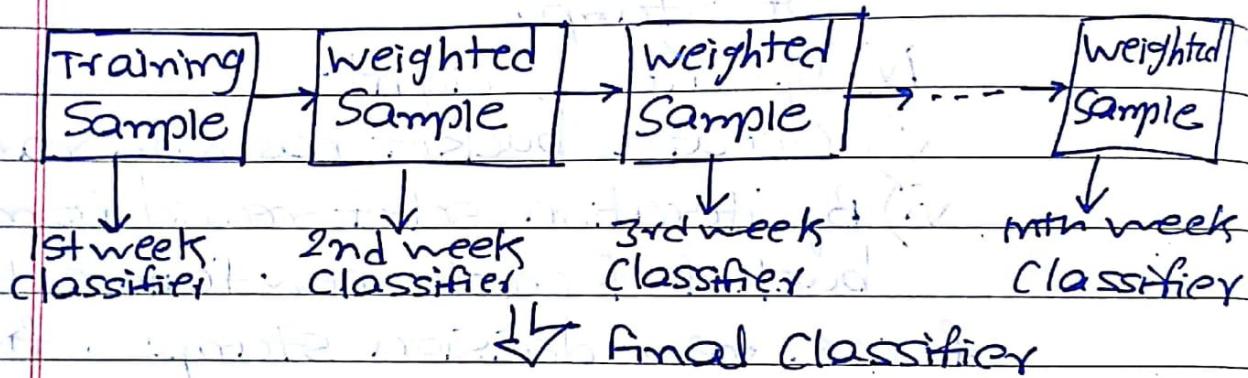
Less error means very less record are incorrectly classified.

So at least we are taking noting of all stumps on test dataset.

So in boosting we are combining weak learner & creating strong learner.

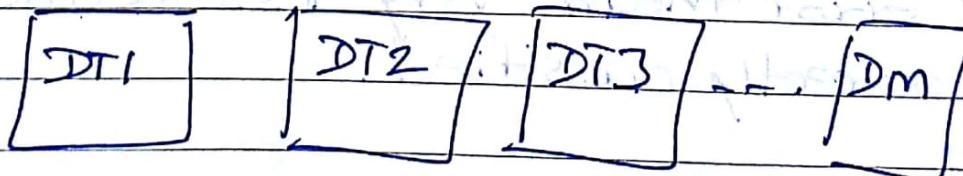
* XGboost

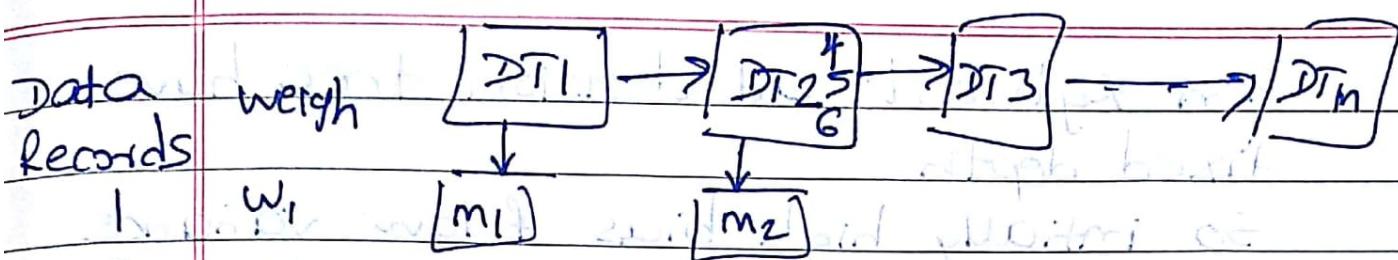
* Extreme Gradient Boosting



→ In XGboost decision tree are created sequentially not parallelly.

→ Sequential Ensemble Technique





1. weigh
2. w_1
3. w_2
4. w_3 Record 1 to 7
5. w_4 Suppose Record 1 to 7
6. w_5 4, 5, 6
7. w_6 are missclassified.
8. w_7 Res weight for the same

- 1) Create a sequential decision tree.
- 2) Assign a weight to each record, but weight is based on probability
- 3) weight decide what is the probability of the selection of record

Initial this weight are equal, that probability of all data record are equal to select for decision tree.

- 4) As per misclassification update weight, increase weight & apply that dataset to next decision tree & based on voting select final classifier output.

In xgboost all decision tree have fixed depth

So initially high bias & low variance which get converted into low bias & low variance. By combining weak learner & get final classifier with good accuracy.

Where as in decision tree, we get low bias & high variance.

Learning from scratch
and I am about 100% done with it.
I am doing it based on boost
which is not a tree which is why
it is more biased to one side but it
is still learning and updating
itself and it is doing it by finding
the variance and picking out nodes &
making cuts based on it.
The other thing is that
it starts with the strongest
bias which is high variance
and then it goes down.