

①

Ensemble Techniques (Classification + Regression)

(We ensemble or club various models together to make a prediction)

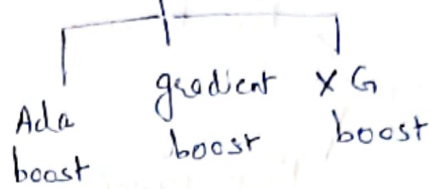
Bagging



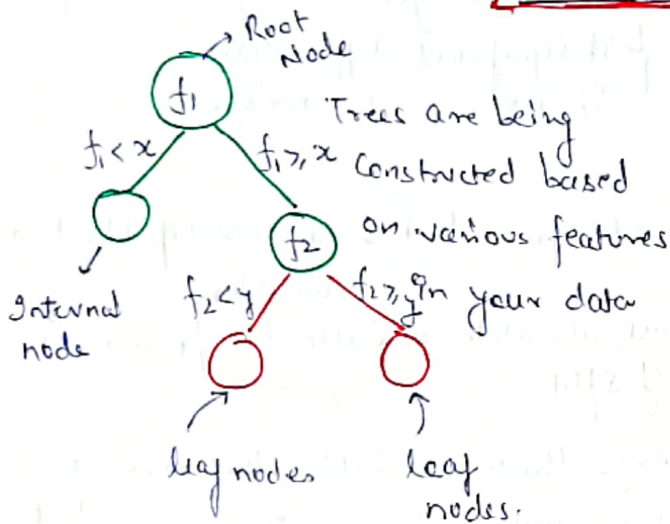
Random Forest

≈ (Decision Tree)

Boosting

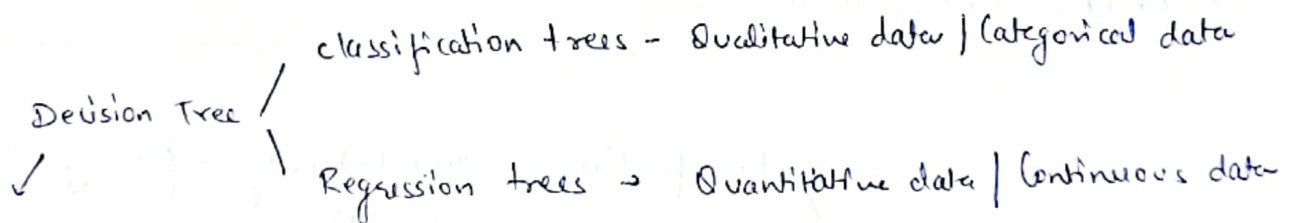


Decision Tree



Lets data

f_1	f_2	f_3	O/P (binary)
			Y
			N



There are very popular methods owing to interoperability of this technique across various divisions.

Simple decision trees are preferred over complex decision trees to prevent overfitting.

- * Multiple way decision tree can be better sometimes than a binary decision tree because it increases the gain in information manifolds.
- * As the number of splits increase in decision trees their complexity rises - Pruning.

In order to prevent overfitting in decision trees we adopt a method called pruning which actually curtails the branches of the tree that fit data too specifically.

(2)

Pruning

- Pre-Pruning** - stop growing DT branches when information become unreliable
- Post Pruning**: first you grow a full-fledged DT and then start pruning irrelevant branches.

Decision Tree Algorithms

CHAID
(Chi-squared Automatic Interaction Detection)

Regression:-

F-test \rightarrow target means/nodes/leaves.
 \rightarrow if significant diff: split
 \rightarrow if not: merge

Classification:- Chi square test \rightarrow relationship b/w two variables

most reliable variable to o/p - node of split
 Merge those variables that are not very significantly related to the output

Short comings:-

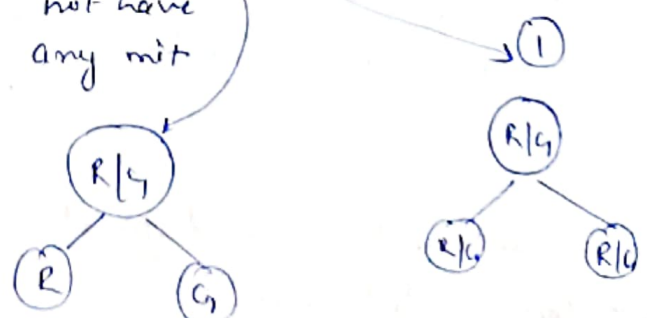
- ① Very horizontal DTs
- ② does not handle missing values handled as a different class.
- ③ no pruning option available

classification:- Gini Impurity - purity of split

$(0-1)$

split does not have any mix

\rightarrow the split has mix of both classes in each node



CART

(binary classification/
Regression Trees)

Regression

Capture more information
 \rightarrow Least Square Deviation (LSD)
 Variance reduction, minimizes
 Sum of sq distances, deviations
 Predicted - observed, sq residual \downarrow

CART produces a sequence of DTs, each of which is a candidate for ③ "optimal tree". This optimal tree is identified by evaluating the performance of every tree through testing or performing cross-validation

does not use any internal performance measure

ID3

used for classification tasks, not very effective with regressions)

Information Gain

→ favouritism towards attributes with more ^{entire} decrease in level of randomness in a set of data

how much information, a feature gives us about a class. attribute with highest information gain will split first.

Entropy: A concept to which information gain is directly related to.

(0-1)
perfectly predictable perfect randomness unpredictable

Information gain $\propto \frac{1}{\text{Entropy}}$

C4.5 → Improvement over ID3

(Regression + Classification Trees)

Gain Ratio: attribute with maximum gain ratio is selected as the splitting attribute
reduces bias in DT with huge amount of branches by taking into considerations number & size of branches while choosing an attribute

Pruning

+

technique

Windowing

→

DT is trained first on a batch of data from training data and rest of the cases are used as a test measure for this DT. If it fits well process stops. Save memory & Computations.

Limitations of Decision Trees

(4)

- ① High Variance - small change in data can result in various sets of splits.
- ② High Bias :- If some classes dominate over others. (Problem with unbalanced datasets)
- ③ Greedy :- locally optimal rather than globally optimal
- ④ Regression :- boundaries.

All of these problems of decision trees can be resolved by ensemble techniques

- Bagging
- Boosting

Mathematics behind DTs

Chi-square test: Hypothesis test when one wants to determine the relationship b/w two categorical variables.

Gender	Preferred NP	frequency of TV	Highest education
1. M	1. Washington P.	1. daily	1. without grad
2. F	2. Hindu	2. several times/week	2. college
	3. USA today	3. more rarely	3. bachelor's

Is there a relationship b/w
or a correlation

↳ chi-square test

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

chi square value

gender < preferred NP. | frequ & highest education

O_i → observed value

E_i → expected value

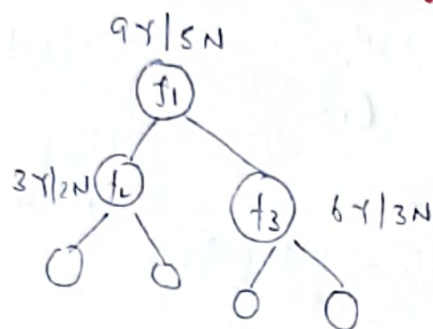
if chi square value < critical chi square value
null hypothesis is ^{not} rejected

From table

Hence there is no relationship b/w these variables.

(5)

Decision Tree Entropy : Measures the purity of the split



$$H(S) = -P_{(+)} \log_2(P_{+}) - P_{(-)} \log_2(P_{-})$$

$$P_{+}/P_{-} = \% \text{ of +ve class} / \% \text{ of -ve class}$$

$$H(S) = -3/5 \log_2(3/5) - 2/5 \log_2(2/5) = 0.78 \text{ bits.}$$

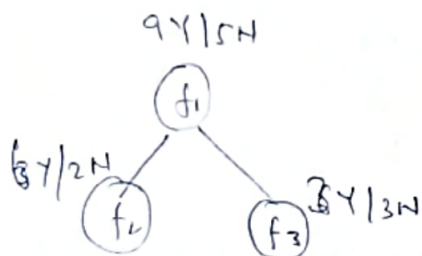
Entropy 0-1

1 - worst split 0 -> best split

Decision Tree Information Gain

$$\text{Gain}(S, A) = H(S) - \sum_{v \in \text{val}} \frac{|S_v|}{|S|} H(S_v)$$

\uparrow subset after splitting
 \uparrow subset entropy



$$H(S) = H(f_1) = -\frac{9}{14} \log_2(9/14) - \frac{5}{14} \log_2(5/14) = 0.94$$

$$\boxed{H(f_2) = 0.91 \quad H(f_3) = 1} \quad H(S_v)$$

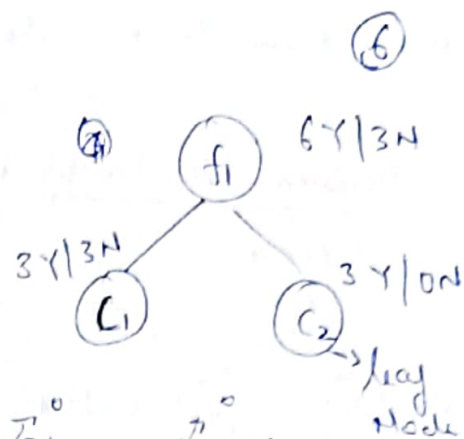
$$\text{Gain}(S, f_1) = H(S) - \frac{8}{14} H(f_2) - \frac{6}{14} H(f_3) = 0.049.$$

Gini Impurity

$$GI = 1 - \sum_{i=1}^n (P_i)^2 = 1 - [(P_{+})^2 + (P_{-})^2]$$

Entropy $H(S) = -p_1 \log_2 p_1 - p_2 \log_2 p_2$

f_1	f_2	f_3	O/P
c_1	d_1		Y
c_2	d_2		N

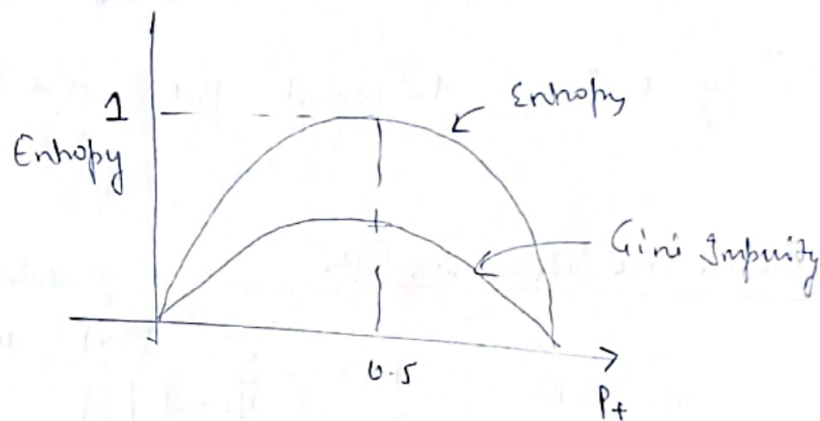


$$H(S) = -\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3}$$

$$H(c_1) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} = 0 \quad (\text{Pure split})$$

$$= 1$$

$$H(c_2) = -\frac{3}{3} \log_2 \frac{3}{3}$$



$$GI(c_1) = 1 - [(1/2)^2 + (1/2)^2]$$

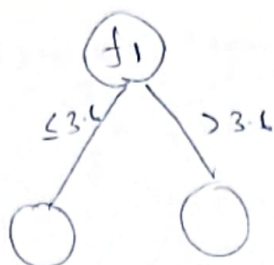
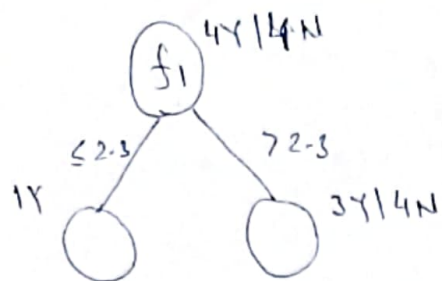
$$= 1/2$$

Gini impurity is taken as a parameter by random forest & XG boost because it is computationally efficient

How decision tree split Numerical Variable

f_1	O/P
2.3	Y
3.6	N
4	Y
5.2	Y
6.7	Y
8.9	Y
10.1	Y

- Sort all the values
- Will consider some threshold values lets say
 $2.3 \quad x_i \leq 2.3 \quad \hookrightarrow \text{branch 1}$
 $x_i > 2.3 \quad \hookrightarrow \text{branch 2}$



(7)

Like this we go with each and every feature

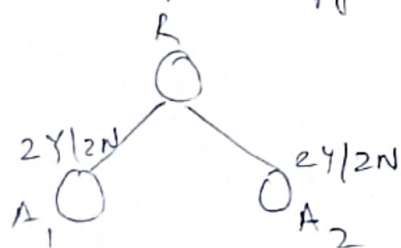
We will calculate the best entropy and information gains.

Disadvantage is time complexity for millions of rows in performing this operations

Gain Ratio

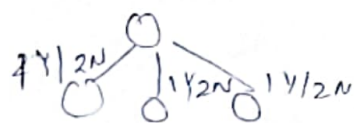
$$\text{Split Entropy } (R, A) = - \sum_{i=1}^P \frac{|R_i|}{|R|} \log \left(\frac{|R_i|}{|R|} \right)$$

Split Entropy increases with the number of divisions. increases



$$= -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2}$$

$$= 1$$



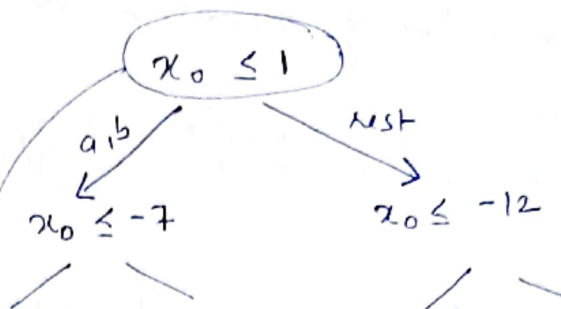
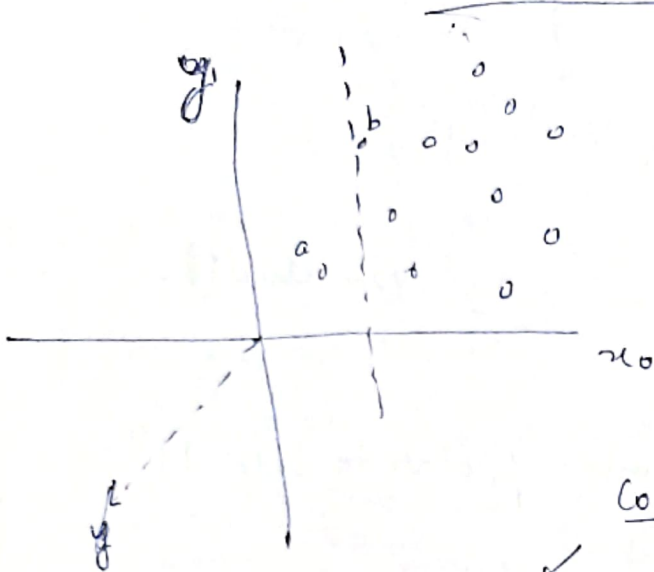
$$= -\frac{1}{3} \log \frac{1}{3} - \frac{1}{3} \log \frac{1}{3} - \frac{1}{3} \log \frac{1}{3}$$

$$= \log_2 3 = 1.6$$

$$\text{Gain Ratio} = \frac{\text{Gain } (R, A)}{\text{Split } (R, A)} \rightarrow \text{Prevent the number of Branches.}$$

Its same as MNP with loss fn as MAE / MSE.

Decision Tree Regression / Random forest Regression

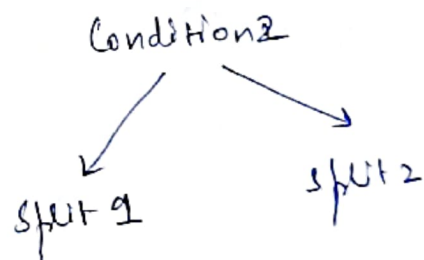
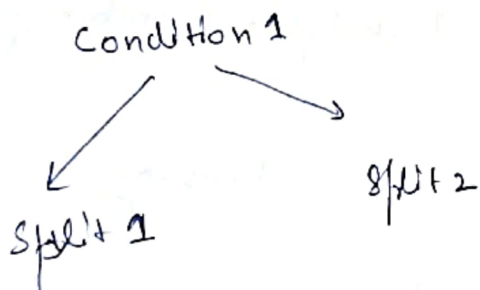


Condition

If this condition is true we move left and if this condition is false we move right

So like this we keep splitting the data until we achieve leaf nodes. It is a sort of clustering if you see.

* Now we have understood the intuition behind decision tree algorithm. Most important bit is to find the condition for splitting at the node



Which of these conditions is better? Condition 1 or Condition 2

Ans- Variance Reduction

In case of regression variance reduction works best like in case of classification we have gini coefficient (6)

$$\text{Var} = \frac{1}{n} \sum (y_i - \hat{y})^2$$

Higher value of variance will mean a poor classification.

$$(\text{Var})_{\text{root}} \geq w_1 (\text{Var})_{\text{split 1}} + w_2 (\text{Var})_{\text{split 2}}$$

$$(\text{Var})_{\text{red}} = (\text{Var})_{\text{parent}} - \sum w_i \text{Var}(\text{child } i)$$

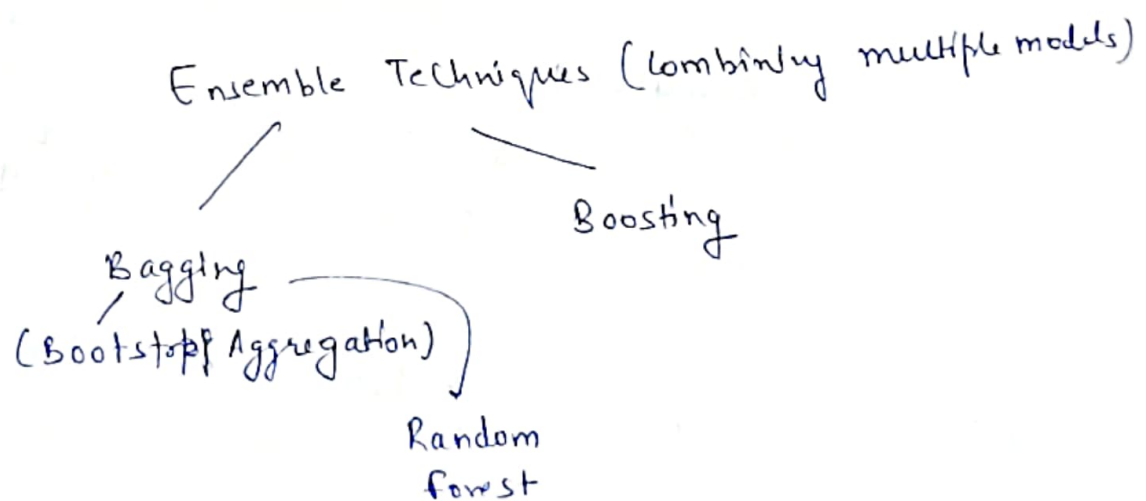
$$w_i = \frac{n_{\text{split } i}}{\text{Total points}} \quad (\text{number of points in split } i)$$

Various split algorithms

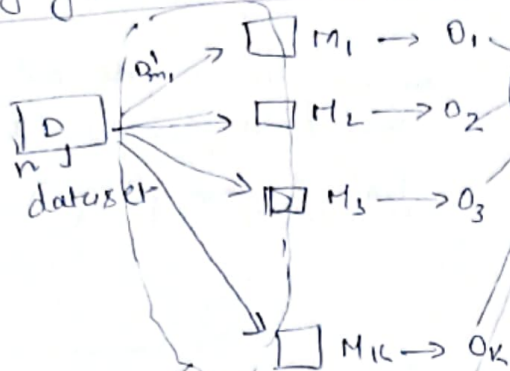
- ID3 - C4.5 - CHAID - MARS

Random forest Regressor (Bagging Technique)

Bagging ??



Bagging | Bootstrap Aggregation



We will create multiple models from a big dataset D and to each model we will provide just a sample of dataset D i.e. D'

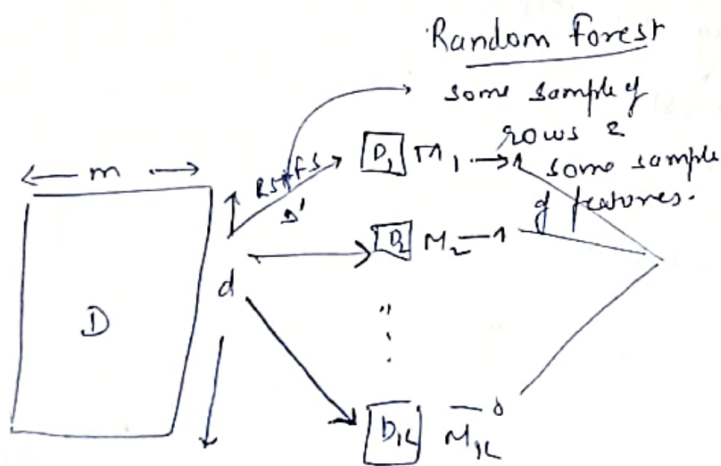
We will again sample the data and provide it to the another model.

m_1 will get trained on D'_1
 m_2 will get trained on D'_2
 \vdots
 m_k " " " " " " D'_k

* $O_1, O_2, O_3, \dots, O_k$ are outputs by different model.

* we will combine all the output and provide the result for test data.

Bootstrap aggregation



$$D' < D$$

There are many decision trees that one creates sampling rows and columns randomly from a big dataset. This sampled data is then passed and used for training different decision trees.

The output of these trees is taken and aggregated to form a common output on a random input.

Base learner in random forest is a decision tree.

Decision Tree

Low bias and high variance

when we create a decision tree to its complete depth it leads to overfitting

✓
Decision tree in complete depth will have a high variance

Random forest overcome this limitation of overfitting of decision tree. As we provide random rows and random features to decision tree

low variance

It is owing to this fact that random forest works very well in most of the cases that we use.

In regression \rightarrow we take mean or median of the decision tree output that we have

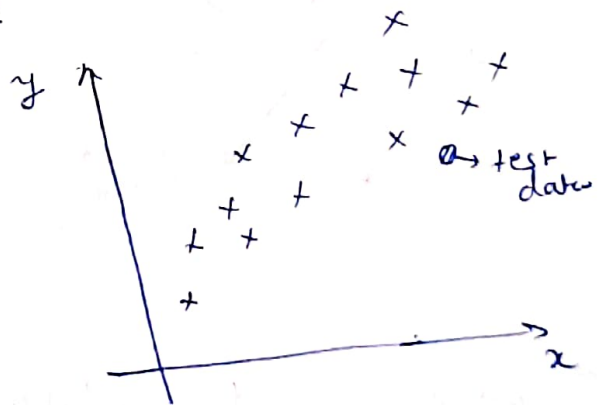
*Hyperparameter \rightarrow how many decision tree we have to use using Random forest

K Nearest Neighbour

classification

Regression

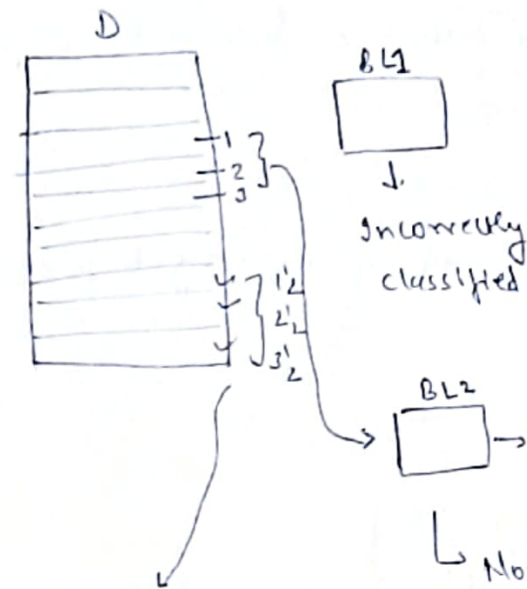
KNN Regression



$K=5$, hyperparameter
i.e. we will look for 5 nearest points to the test point
average of all the values of 5 data points
value = average of nearest datapoints.

Boosting

(11)



① we pass entire data to the base learner and then we see that which of the data has been incorrectly classified.

② Let's say record 1, 2, 3 are incorrectly classified.

These records will now pass to Base learner 3



Now let's say base learner 2 gives wrong values for other three records

This cycle will keep going on in the same manner. This is a boosting technique

Sample dataset

Adaboost → Base learner - Decision Tree

f_1	f_2	f_3	O/P	Sample weight
-	-	-	-	$1/7$
-	-	-	-	$1/7$
-	-	-	-	$1/7$
-	-	-	-	$1/7$
-	-	-	-	$1/7$
-	-	-	-	$1/7$

sample total error = $1/7$

Cycle 1 formula for Sample weight

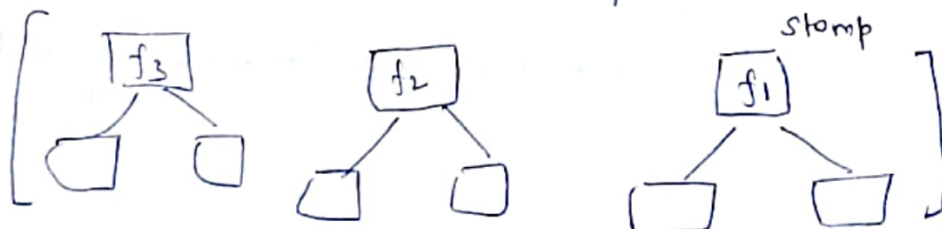
$$w = 1/n$$

$$= 1/7$$

all records are basically assign the same weight

* Decision Tree are base learners and are created with just one depth. These are called as stumps

We create Stump for each feature with a depth 1



12 We will select base learner as a decision tree which has least entropy
Now let's say it correctly classifies 3 records & incorrectly classifies 4 records.

Step 2 Calculate the total error - just sum up the weights for the error rows.

Step 3

$$\text{Performance}_{\text{stump}} = \frac{1}{2} \log_e \left[\frac{1 - TE}{TE} \right]$$

$$TE \rightarrow \text{Total error} = \frac{1}{2} \log_e [6] = 0.896.$$

Step 4

update the weights \rightarrow Increase the weight for wrongly classified records and decrease the weights for correctly classified ones

Update weight

$$\begin{aligned} \text{Incorrectly classified record} \quad \text{New weight} &= \text{Old weight} \times e^{\text{Performance}_{\text{stump}}} \\ &= \frac{1}{7} \times e^{0.859} = 0.359 \end{aligned}$$

$$\begin{aligned} \text{For correctly classified point} \quad \text{New weight} &= e^{-(\text{Performance}_{\text{stump}})} \\ &= e^{-0.895} = 0.05 \end{aligned}$$

Now we will normalize the weights by dividing all these newly obtained weights by sum of these weights.

Create a New dataset for next Base learner

13

Normalised weight

0.07
0.51
0.07
0.07
0.07
0.07

Category

0 - 0.07
0.07 - 0.55
0.55 - 0.65
0.65 - 0.72

New dataset

* Now it will select a random weight from one of these buckets.

Let's say random weight is

0.34

* which bucket it lies?

③

* Select one or two records from this bucket.

* Make a dataset.

Now owing to updated weights there is a more probability of a wrong weight getting selected.

This cycle will continue until it pass through sequential decision tree. We will obtain the sequence of stumps and for a test data we will obtain value from all stumps and will either take mean or median.

* We combine multiple weak learners to obtain a strong learner

Gradient Boost

Exp	Degree	Salary	1/3
2	BE	50K	75
3	Master	70K	75
5	master	80K	75
6	PHD	100K	75

Step 1

Base model - 1 (average)

$$\frac{50+70+80+100}{4} \approx 75K$$

what ever the input is, I will give the output as 75K.

Step 2 Compute pseudo residuals / residuals / error

$y - \hat{y}$

Step 3: Construct the decision tree

Exp	Degree	Salary	\hat{y}	R_1	R_2	R_3
2	BE	50K	75	-25	-23	1
3	MBA	70K	75	-5	-3	1
3	MBA	80K	75	5	3	1
10	PHD	100K	75	25	20	

Output is residual and not all the features

$\{x_i, R\}$

Base learner \rightarrow Decision Tree 1 \rightarrow DT₂ \rightarrow DT₃
 average (x_i, R_2) (x_i, R_3)

How to compute salary

$$75 + R_1 + R_2 + \dots$$

like $75 - R_3 \approx 52$
 which is very near to 50K

Now we need a decision tree with low variance & low bias.

Hence we will introduce a learning rate

This is a problem as my model is overfitting

$$75 + \alpha(R) \rightarrow (0-1)$$

$$F(x) = h_0(x) + \alpha_1 h_1(x) + \alpha_2 h_2(x) + \dots + \alpha_n h_n(x)$$

$$= \sum_{i=1}^n \alpha_i h_i(x) + h_0(x)$$

This residual value will be decrease and one with more residual will be decreasing slowly but other near values residual will eventually come to zero.

Pseudo algorithm (gradient boost)

15

Exp	Degree	Salary	\hat{y}	r_1	r_2
2	BE	50K	60	-10	
3	MBA	70K	60	10	
4	MAST	60K	60	0	

I/P

- ① $\{x_i, y_i\}$: Independent and dependent variables
- ② $L(y, f(x))$: loss function (differentiable)
Regression:- Mean squared error
Classification:- Hinge loss / log loss
- ③ Number of Trees

Pseudo Algorithm

- ① Initialize the model with constant value
$$F_0(x) = \underset{\gamma}{\operatorname{argmin}} \left(\sum_{i=1}^n L(y_i, \gamma) \right)$$

$$\text{Loss} = \sum_{i=1}^n \frac{1}{2} (y_i - \hat{y})^2 \downarrow$$

Find the γ value for which loss f^n is minimal

$$L = \frac{1}{2} (50 - \hat{y})^2 + \frac{1}{2} (70 - \hat{y})^2 + \frac{1}{2} (60 - \hat{y})^2$$

↓ 1st order derivative

$$\frac{\partial L}{\partial \hat{y}} = -(50 - \hat{y}) - (70 - \hat{y}) - (60 - \hat{y}) = 3\hat{y} - 180 \approx 0$$

$\hat{y} = 60$ - average

- ② Iterate the steps from 1 to (M) → number of DT_L

Compute pseudo residuals.

$$r_{im} = \left[\frac{\partial L(y_i, f(x_i))}{\partial (f(x_i))} \right]_{i=1+m}$$

(16) Loss = $\frac{1}{2} (y - \hat{y})^2$

$$\frac{\partial L}{\partial \hat{y}} = (y - \hat{y})(-1)$$

$$\boxed{-\frac{\partial L}{\partial \hat{y}}} = y - \hat{y}$$

$$\hookrightarrow - \left[\frac{\partial L(y, f(x_i))}{\partial f(x_i)} \right]$$

Now this residual error is nothing but the difference of actual & predicted value

Once we get the residual we create a decision Tree where my dependent feature will be the residue and independent feature will be experience and degree

- again fit a base learner. $h_m(x)$

$$i \in \{x_i, x_{im}\}$$

Step 3

$$\gamma_m = \underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^n \left(L(y_i, \overset{\substack{\text{prev model} \\ \text{output}}}{f_{m-1}(x_i)}) + \gamma \right)$$

$$L(y_i, f_{m-1}(x_i)) = \sum_{i=1}^n \frac{1}{2} (y_i - (60 + \hat{y}))^2 - \text{minimize this.}$$

Step 4

$$F_m(x) = F_{m-1}(x) + \alpha (h(x))$$

\hookrightarrow learning rate

XG Boost Regressor

Average salary = 51K

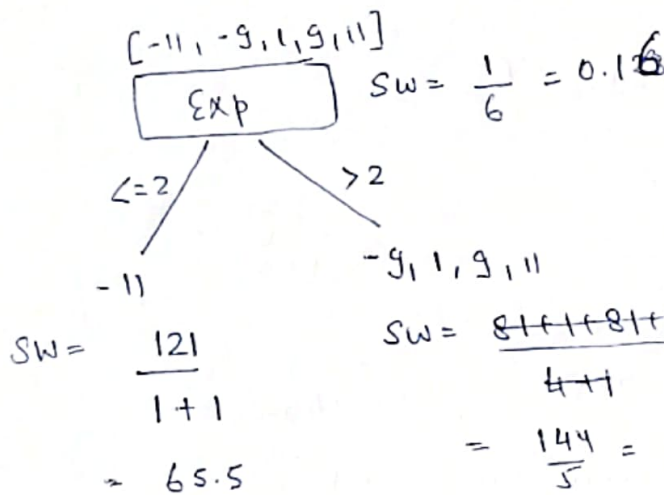
(14)

Base model = 51K

Exp	Cap	Salary	Residual 1	O/P	Res2	→ Another tree
2	Yes	40K	-11K	46	1	
2.5	Yes	42K	-9K	46	1	
3	No	52K	1K	53.5		
4	No	60K	9K	62		
4.5	No	62K	17K	63		

DT (Exp, Cap, Res1)

only binary trees



Step 2: Similarity weights

$$= \frac{\sum (\text{Residual})^2}{\text{No. of Residuals} + \lambda}$$

hyper parameter

$$\lambda \approx 1$$

$$\text{Gain/Information gain} = (LSW + RSW) - \text{Root SW}$$

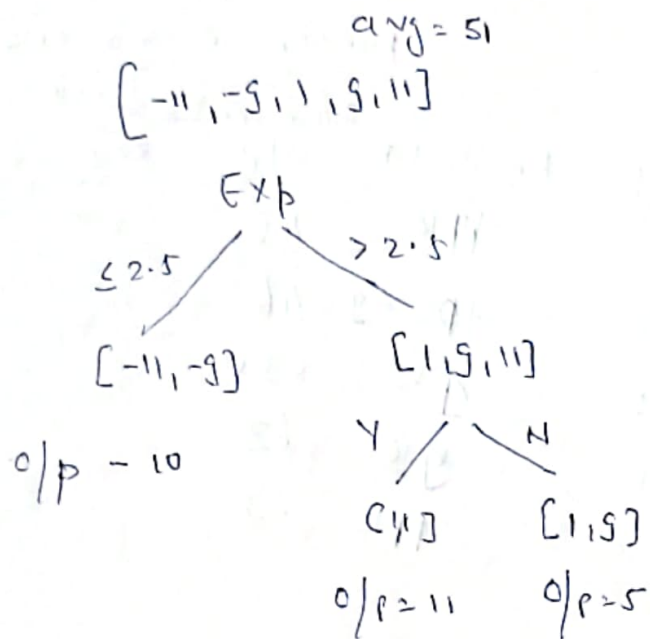
$$= 60.5 + 28.8 - 0.16 = 89.14$$

Now similarly we will make other split and calculate the gain

≤ 2 89.14 / We decide the split

≤ 2.5

We can create any number of decision trees, and select 1 DT and calculate output



Now the value for exp 11

$$= 51 - \alpha(10)$$

Let $\alpha = 0.5$

$$= 46.$$

$$o/p = \text{Base model} + \alpha_1(T_1) + \alpha_2 T_2 + \alpha_3 T_3 \dots + \alpha_n T_n$$

There is one more hyper parameter γ

Let say $\gamma = 150.5$

$$\text{a gain} - \gamma < 0$$

prune this tree means

↳ Cut it → prevent the overfit

if positive don't prune this.

