

Agenda:

- ① Decision Tree classification
- ② Decision Tree regression
- ③ Practical implementation
- ④ Ensemble technique

first we check for Categorical data, then for Continuous data

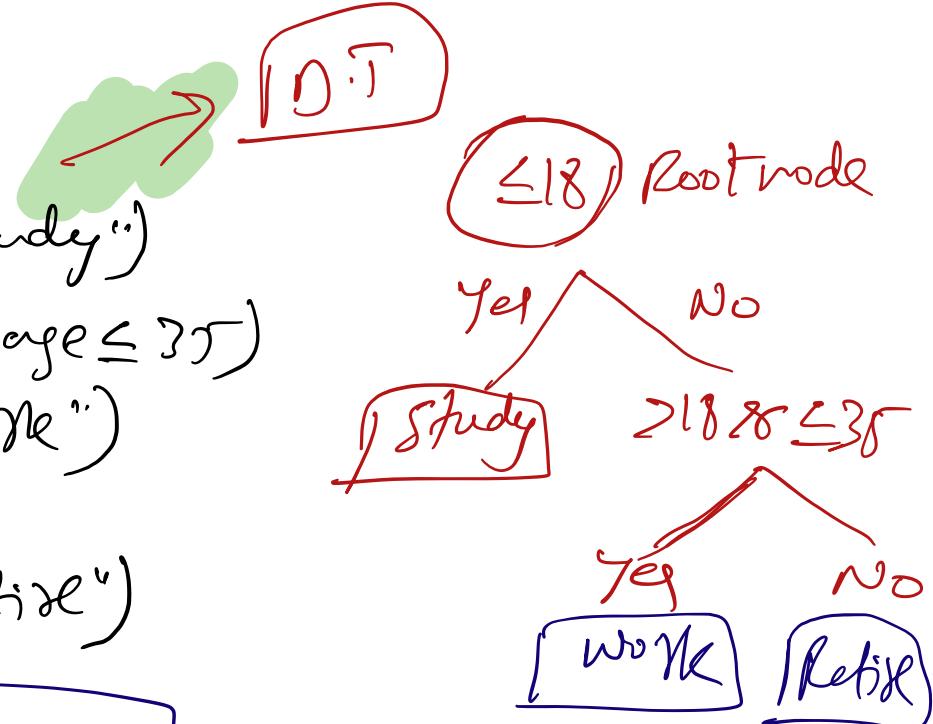
→ we can solve both classification and Regression Problems.

ex:

```

if (age <= 18)
    print ("study")
elif (age) > 18 & age <= 35)
    print ("work")
else:
    print ("Retire")

```

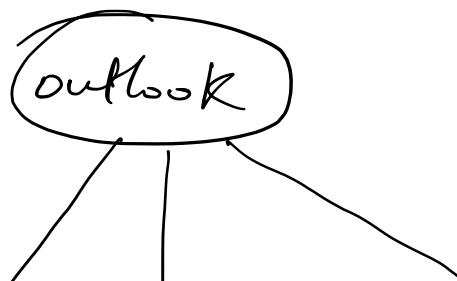


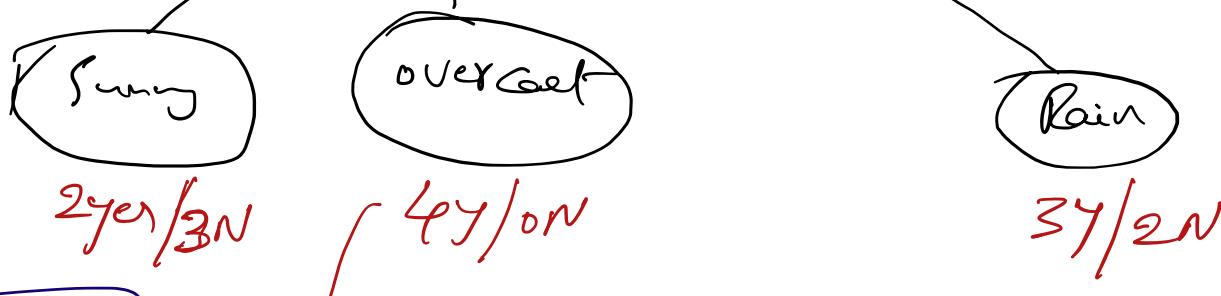
note:

Nested if-dse \Rightarrow DT

ex: dataset | Day | outlook | temp | humidity | wind | playTennis
yes/no

Classification:





Cut overcast

Cut, there is no confusion
it is clear they will
play.

→ we split that particular node until we reach a pure node.

→ this process goes on.

(AKA pure leaf node)

No need to split it even more.

We use two methods for that called:-

① Entropy
② Gini Impurity

② How the features are selected?

Note: ① Purity

→ to find whether the split is pure (Y) or not (N)

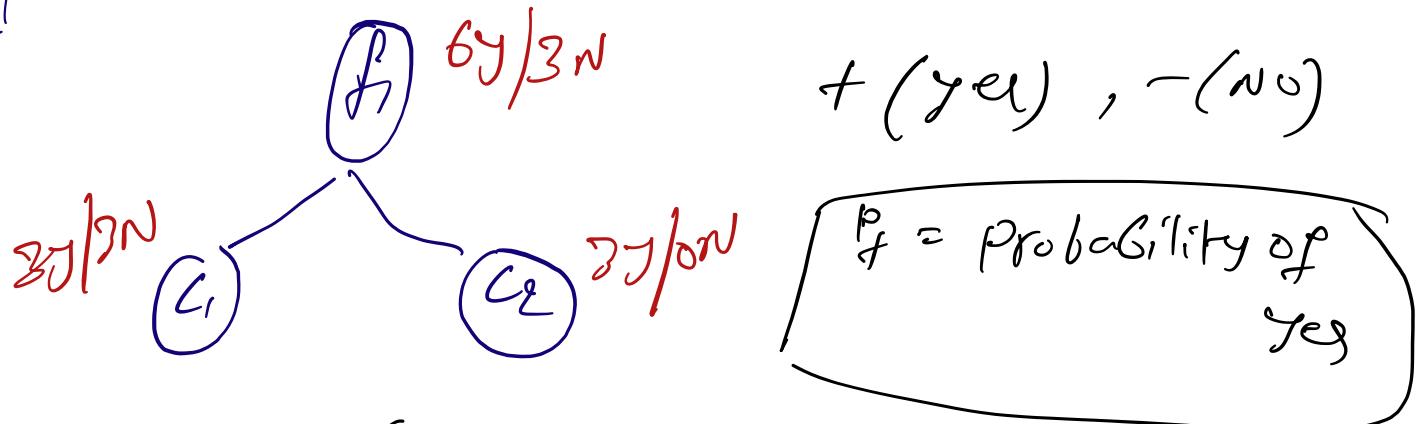
→ Is the the node pure/leaf node or not

⇒ Information Gain

Entropy:

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

ex:

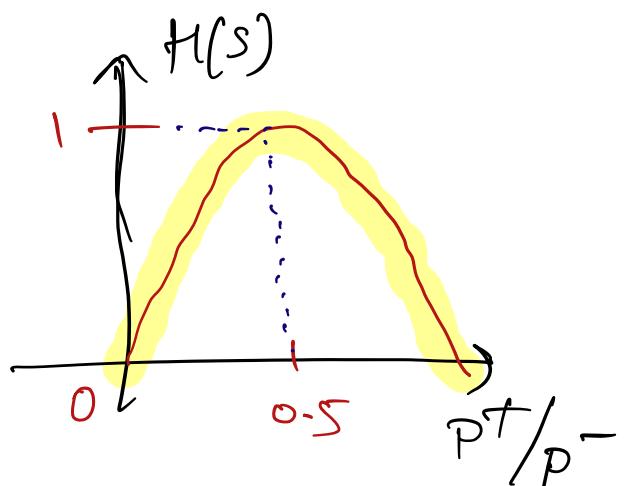


$$P_+ = -\left(\frac{2}{3}\right) \log_2 \left(\frac{2}{3}\right) - \left(\frac{0}{3}\right) \log_2 \left(0/3\right) \rightarrow (\text{pure split})$$

$$C_2 = -1 \log_2 1 = 0 \Rightarrow \text{hence, it is a pure split.}$$

P_- = probability of No

$$C_1 = -\frac{3}{6} \log_2 \frac{3/6}{} - \frac{3}{6} \log_2 \frac{3/6}{} = 1 \quad (\text{impure split})$$



Note: Entropy is always b/w 0 and 1

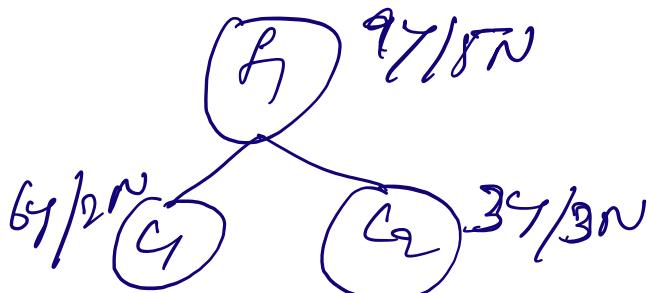
Information Gain :-

Root node Entropy

that particular node entropy

$$\text{Gain}(S, f_1) = H(S) - \sum_{V \in \text{Val}} \frac{|S_V|}{|S|} H(S_V)$$

Ex:-



$H(S) \Rightarrow$ Entropy of the Root node

$$\begin{aligned} H(S) &= -P_+ \log_2 P_+ - (P_-) \log_2 (P_-) \\ &= -\left(\frac{9}{14}\right) \log_2 \left(\frac{9}{14}\right) - \frac{5}{14} \log_2 \left(\frac{5}{14}\right) \end{aligned}$$

$$H(S_V) = -\frac{6}{8} \log_2 \left(\frac{6}{8}\right) - \frac{2}{8} \log_2 \left(\frac{2}{8}\right)$$

$$H(C_1) \approx 0.81$$

$$H(C_2) = 1$$

$$\begin{aligned} \text{Gain}(S, f_1) &= 0.94 - \left[\frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right] \\ &= 0.041 \end{aligned}$$

→ if we get the gain for $f_2 = 0.051$
 then what to use?

$\text{gain}(S, f_2) \gg \text{gain}(S, f_1)$

(we chose this)

Note: T-train should be more to be selected

Gain Impurity:

range

0 to 0.5

$$H(S) = 1 - \sum_{i=1}^n (P_i)^2$$

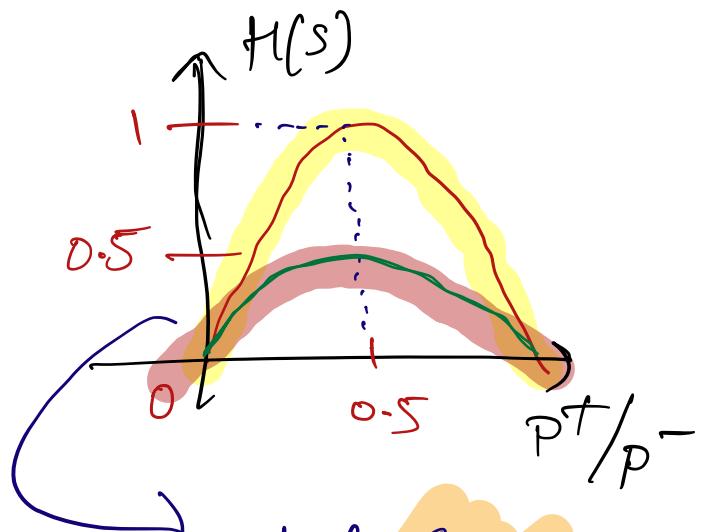
$$= 1 - [(P_+)^2 + (P_-)^2]$$

ext

$D_{27/2N}$

we got entropy = 1
 for same above

Now,



$$h_i = 1 - \left(\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right)$$

$$= 1 - \left(\frac{1}{2}\right) = 0.5$$

Graph for Gini Impurity

When to use wt?

→ Decision tree itself have very bad time complexity
(\propto execution time)

→ Entropy takes more time to execute than Gini Impurity

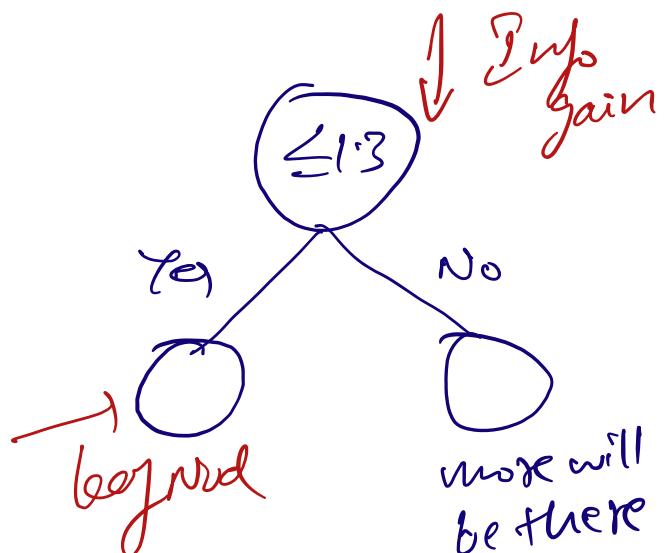
Note!

⇒ If we have 100's of features, then we use Gini Impurity

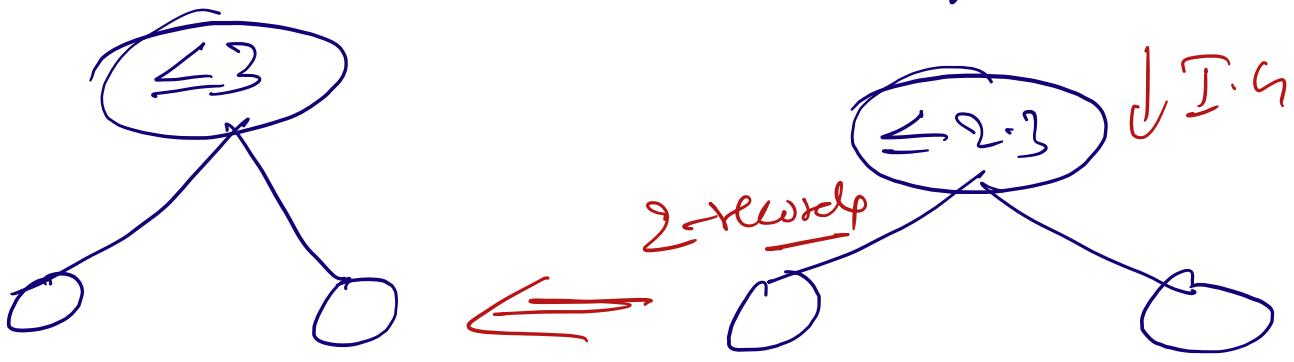
⇒ If we have few features, then we use Entropy.

For Continuous data

ext	f _i	op	$\Rightarrow f_i$
2.3			1.3
1.3			2.3
4			3
5			4
7			5
3			7



next it will take

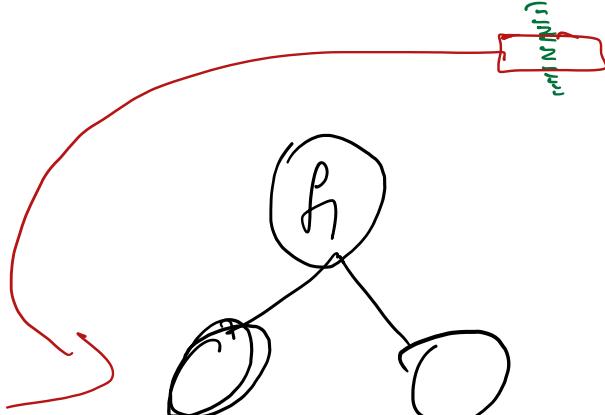


Note: In this way they calculate for every feature making them root and then → then we take the height value node and then split them further.

Decision Tree Regression



Continues



→ left P_1 mean value is calculated

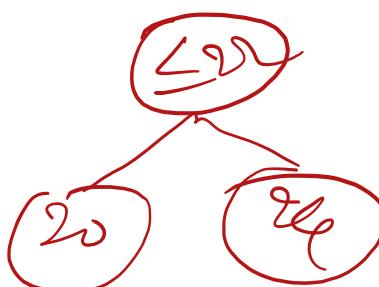
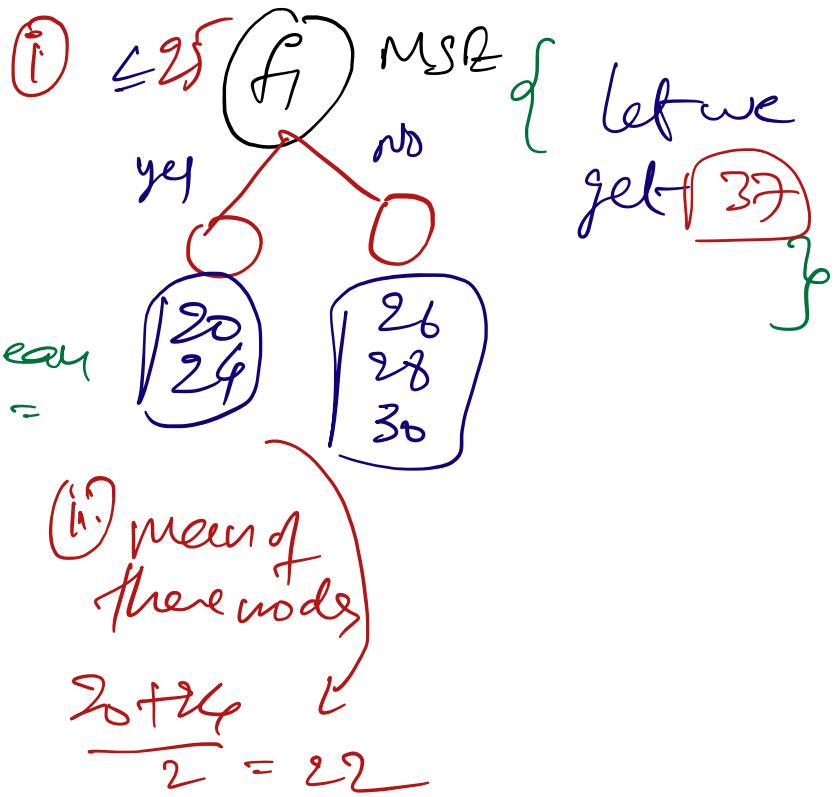
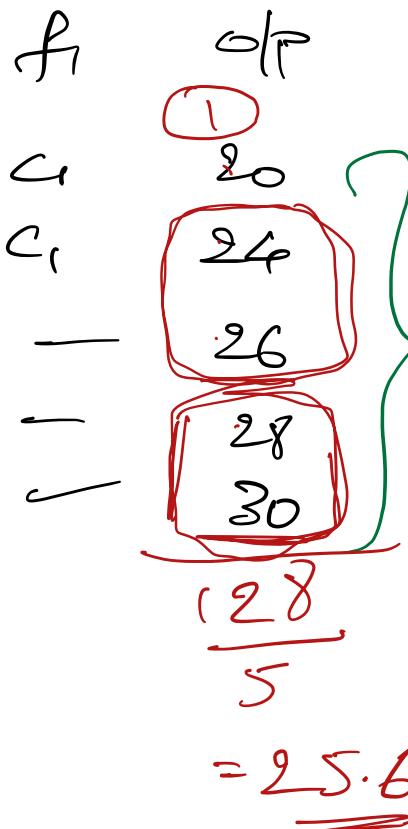
→ Compute MSE
based on that, we split.

$$\frac{1}{2m} \sum_{i=1}^m (y_i - \bar{y})^2$$

after take values and keep in a node and then again pick it for MSE again.

and this process continues. As the M.R. value reduces it means we are nearing leaf node. and on the leaf node, whatever mean we have it is our O.P.

Ex:-



Note: Decision Tree leads to overfitting, to deal it

- ① post pruning
- ② pre pruning

Ex:- ODDS / EVO So, it roughly say that it gives node

so, we don't do more pruning and stop here, called post pruning.

pre pruning:

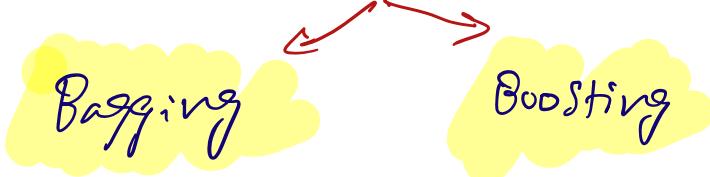
- it is calculated by hyperparameters
- it tells abt max.depth, max.leaf



Agenda:

- (1) Ensemble Techniques
 - Bagging
 - Boosting
- (2) Random Forest
- (3) AdaBoost
- (4) XG Boost

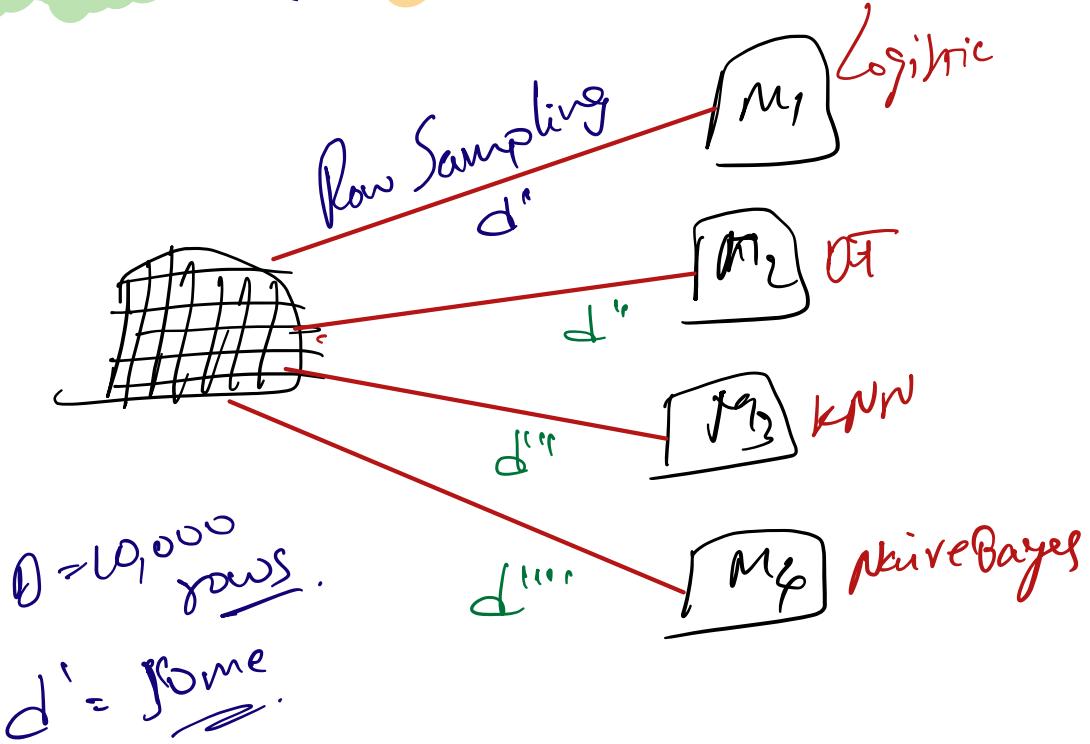
Ensemble Technique: we were discussing only single algo at a time, now we see how to use multiple algos together.



Bagging!

(for classification)

$D' \ll D$



Note: when doing row sampling some of the rows might be repeated in other algos

How the inferencing will happen in test data:



give $o/p = 0$ } binary classification



give $o/p = 1$



give $o/p = 1$



give $o/p = 1$

} all these o/p will get aggregated and the result will be acc to majority voting.

{ Bootstrap Aggregator }

Note:

What if we have equal no. of 1's and 0's,
there is no chance for them to be tie... .

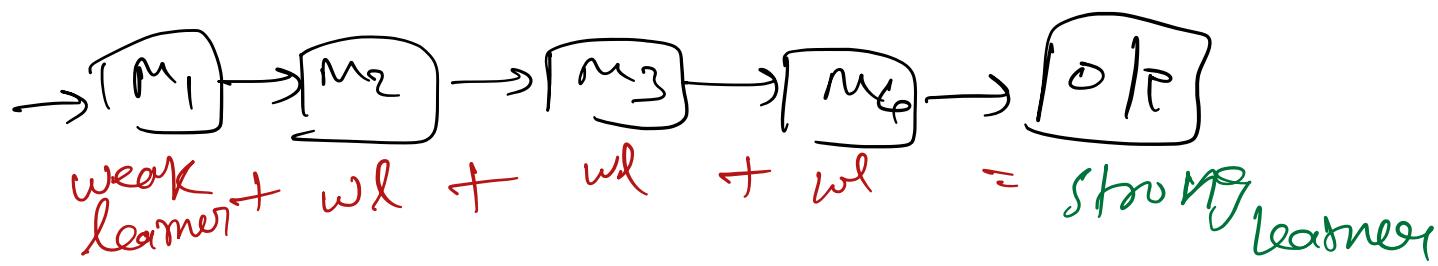
(for Regression)

$$\begin{aligned}M_1 &= 120 \\M_2 &= 140 \\M_3 &= 122 \\M_4 &= 148\end{aligned}$$

} o/p

We take mean of these values,
that will be the o/p of the
model.

Boosting: It is a sequential combination of
weak learning models.



Bagging

- (1) Random Forest Classifier
- (2) Random Forest Regression

Boosting

- (1) AdaBoost
- (2) Gradient
- (3) Xgboost

$\hookrightarrow \underline{=}$

In DT we have problem of overfitting

$\left\{ \begin{array}{l} \text{low Bias} \\ \text{High Variance} \end{array} \right\} \Rightarrow$ we convert
 it to $\left\{ \begin{array}{l} \text{low Bias} \\ \text{low Variance} \end{array} \right\}$



- Normalization / Standardization is not required in DT
 \Rightarrow It is required in KNN
- Random is not impacted by the outliers
 KNN is impacted.

Adaboost:

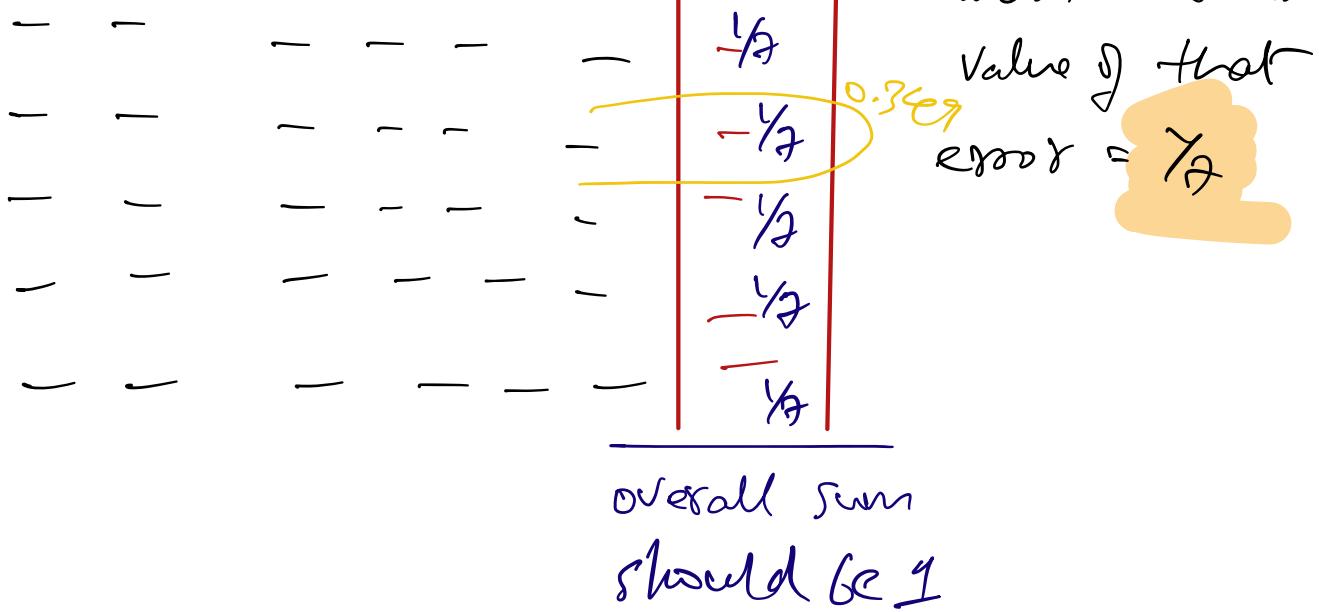
f_1	f_2	f_3	f_4
—	—	—	—
—	—	—	—

O/P	Yes	No
WT	X ₁	X ₂

(1.)

note:

0.05 if there is an error, then the



② information gain



depth will be
of only one level
called **Stumps**

[weak learner]

① Total error

Performance of Stump :-

$$② = \frac{1}{2} \log_e \left(\frac{1 - TE}{TE} \right)$$

$$\text{ex:- } TE = \frac{1}{7} \Rightarrow P = 0.895$$

'3' main
things
In AdaBoost

③ update wts

Note: The corr record wts will be reduced, and the error records wt will be increased.

New sample wt =

Correct Record:

$$(\text{wt}) \times e^{-ps}$$

→ The wrong record wts will then go to the next weak learner

$$\Rightarrow \frac{1}{2} \times e^{-0.895} = 0.05$$

In Correct wts:-

$$(\text{wt}) \times e^{ps} \Rightarrow \frac{1}{2} \times e^{0.895} = 0.349$$

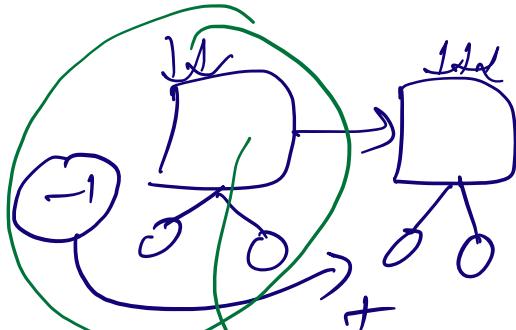
⇒ If we find the total sum of these new wts we find it is not equal to 1

- ① For this reason, we use **Normalization wts**:
- ② we divide each new wt with the new wts summation.

Ex:-

$$0.05 / 0.649 = 0.07$$

②



this bucket randomly creates num b/w 0 to 1

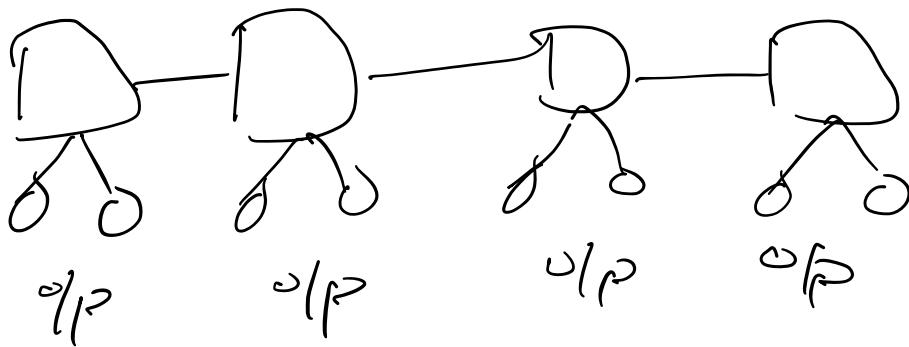
- first model we create, we got an error
- ↳ now along with that error, we pass some more rows into a new stump with the help of first-stump.
- ↳ for this we have a concept **Bucket**

Buckets

- $[0 \text{ to } 0.07]$ → this record is getting added to next row
 - $(0.07 \text{ to } 0.14]$
 - $(0.14 \text{ to } 0.21]$
 - $(0.21 \text{ to } 0.247)$
 - $(0.247 \text{ to } 0.251)$
- 
- 

probability of that number being in any of these buckets, it will be more for the **cover bucket**, so along with other buckets, this will also get passed.

After creating stamps of various DT's:



ex: 0 1 1 1



again we take majority voting
(for classifier)

$$= \frac{1}{4} (0/p)$$

Avg for regressor:

Black box (Vs) White box:

Linear Regression \rightarrow white Box

Random Forest \rightarrow Black Box

DT \rightarrow white Box

ANN \rightarrow Black Box

Interview
Qn
=

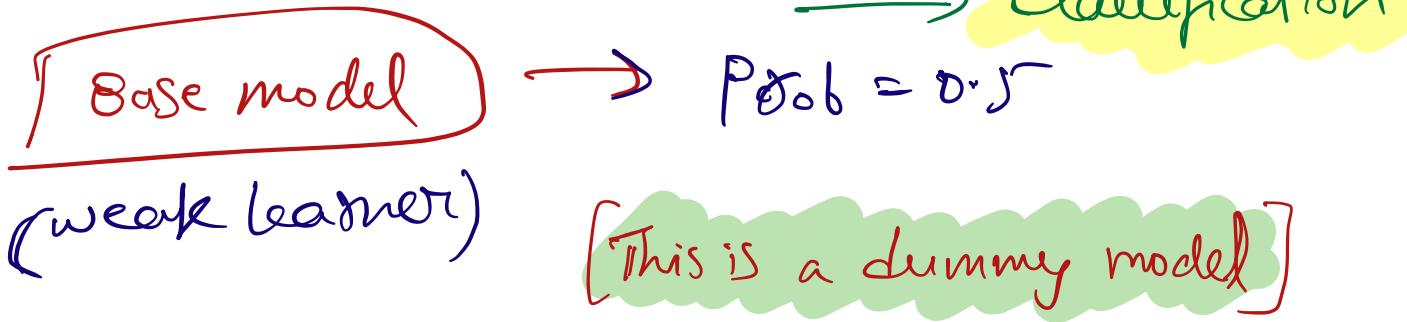
Agenda:

- (1) Xgboost classifier
- (2) Xgboost Regressor
- (3) SVM
- (4) SVR

Xgboost Classifier: [Extreme Gradient Boosting]

ex: {Dataset}	<u>Salary</u>	<u>Credit</u>	<u>Approval</u>
	≤ 50	Bad	0
	≤ 50	Good	1
	≤ 50	A	1
	> 50	B	0
	> 50	C	1
	$> 50k$	Normal.	1

- This Xgboost can be used for both **Classification**, and **Regression** problems
- Internally we use **DJ** in this Xgboost
- When we use Xgboost, we first create a **Base model** with probability = 0.5



→ Then we create a field in the original table called as **Residual Field**
 (probably a new column)

Salary	Credit	Approval	Residual	Approval - base model Prob
<=50	Bad	0	→ (-0.5)	
<=50	Good	1	→ 0.5	
<=50	A	1	→ 0.5	
>50	B	0	→ (-0.5)	
>50	G	1	→ 0.5	
>50k	Nor	1	→ 0.5	
<=50k	Normal	1	→ 0.5	
<50k	Normal	0	→ (-0.5)	

Steps: After creating a base model

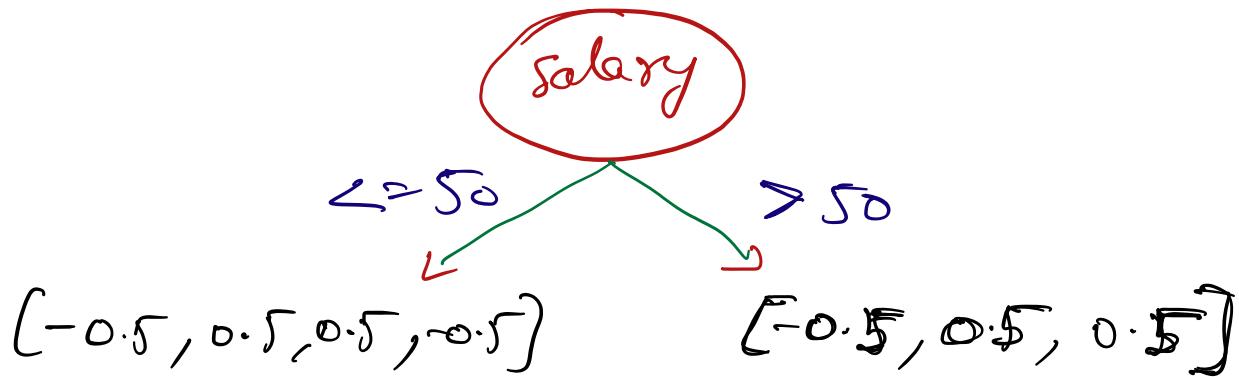
- ① Create a **Binary Decision Tree** using the features
- ② Calculate the **similarity weight**

→ $\frac{\sum (\text{Residual})^2}{\text{Total}}$

③ Calculate the Information Gain

we use these steps to create a Xgboost classifier.

→ we always create a binary decision tree



(Step ① is done)



↳ Now for this, we get

Similarity wt

$$= 0.33$$

↳ ②

$$= \frac{(-0.5 + 0.5 + 0.5 - 0.5)^2}{0}$$

$$\begin{aligned} & 0.5(1-0.5) + 0.5(1-0.5) \\ & + 0.5(1-0.5) + 0.5(1-0.5) \end{aligned}$$

Similarity wt = 0

we generally initialize $\lambda = 0$

1

→ if we calculate similarity wt for root node
↳ take all residuals

$$= 0.142$$

2 →

3

(Step ② is done)



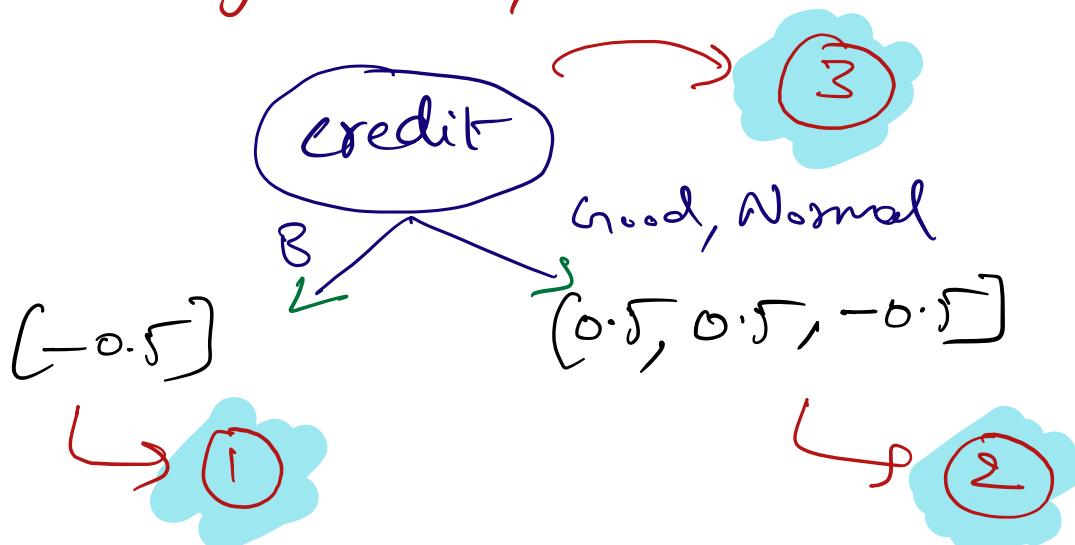
$$\text{Information Gain} = D + 0.33 - 0.14 \\ = 0.19$$

note: we did

$$(1+2-3)$$

→ after we get the information gain, we use it to do further split and select the exact feature, which helps us in finding the proper root node for the model.

→ after salary, we split credit



find similarity wt :-

$$\begin{array}{l} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \end{array} \quad \begin{array}{l} 1 \\ 0.33 \\ 0 \end{array} \quad \left. \begin{array}{c} \\ \\ \end{array} \right\} \quad 1 + 0.33 - 0 = 1.33$$

↳ this "zero" came from $\textcircled{1}$ of salary tree
with $L = 50$

Note:

and this process goes on.

Inference :-

- ① The first record will go to the base model and the probability is 0.5
- ② But, this is an assumed probability,
how to find the actual probability???
- ③ we use log s for that.

$\log(P/(1-P))$ this is applied only for the base model.

$$\hookrightarrow \text{in this case } \log \frac{0.5}{1-0.5} = \log 1 = 0$$

now, when a record goes into this base model

we get the probability =

for 1st row = $[o + \alpha(1)]$



$$\sigma [o + \alpha(1)]$$

actual prob

the first row

goes into the
Binary tree and

reaches till Credit, Bad

Similarity wt = 1

Sigmoid
function
(activation
function)

learning
rate

it can be any small value
acc to the learning rate

we choose

(just like gradient descent)

$$\alpha = 0.01$$

based on the
value of α ,

the o/p will be
between 0 to 1

In this case.

Hence;

$$\sigma [o + \alpha_1(D_{11}) + \alpha_2(D_{12}) + \alpha_3(D_{13}) + \dots + \alpha_n(D_{1n})]$$

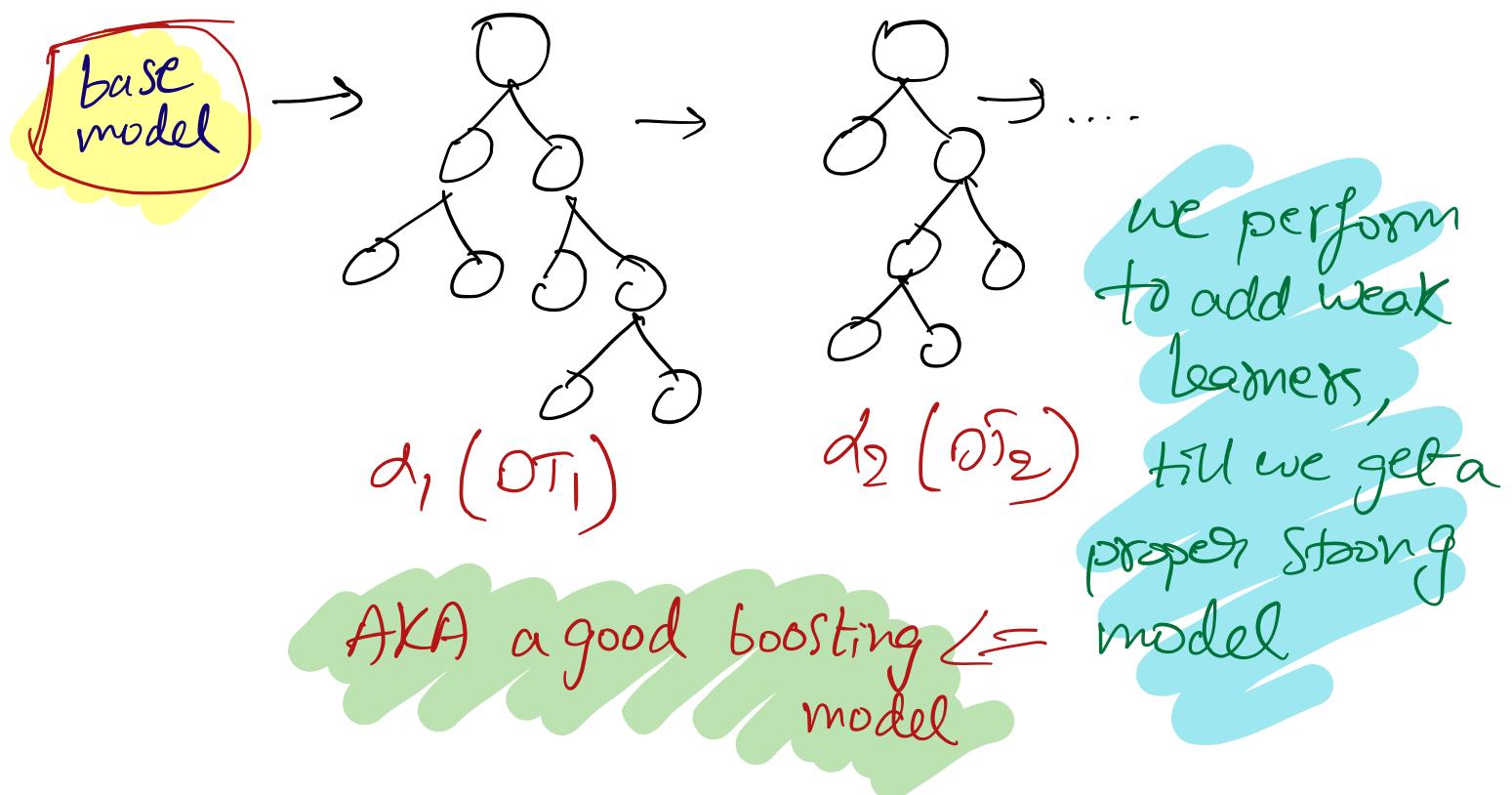
Since,

we are adding all the D_i 's it results us in **Boosting**

Note:

- ① we can split credit even more, but we prefer a node whose **Information Gain** is more
- ② Since, there is lot of math involved, we call Xgboost as a **Black Box model**
- ③ we tend to do **pre-pruning** in their total process

Diagrammatically:-



④ (DT's) get constructed depending on
Independent features

⑤:

we basically setup the hyperparameters [d]
by the help of Cross Validation

Xgboost Regressor:-

Type	Exp	Salary	$\rightarrow O/P Reg$
2	Yes	40K	-11K
2.5	Yes	49K	9K
3	No	52K	1K
4	No	60K	9K
4.5	Yes	62K	11K

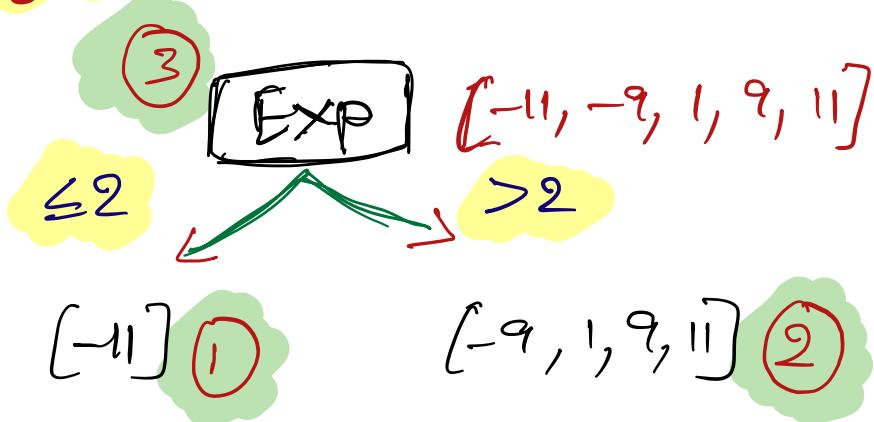
\rightarrow we create base model = avg of O/P's

$$= \frac{\sum \text{Salary}}{6}$$

$$= 51K$$

\rightarrow Residual

\Rightarrow Salary - base model



Similarity wt =

$$\frac{\sum (\text{Residual})^2}{\text{No. of Residuals} + 1}$$

\Rightarrow let $\lambda = 1$

$$\begin{array}{l}
 \textcircled{1} \quad 60.5 \\
 \textcircled{2} \quad 28.8 \\
 \textcircled{3} \quad 0.16
 \end{array}
 \left. \right\} I_{\text{oh}} = 60.5 + 28.8 - 0.16 = 98.13$$

\Rightarrow we got (98.13) acc to the first row with $\boxed{\text{Exp} = 2}$, we will get diff (DT) 's for diff rows

\Rightarrow we pick the split which has the better Information gain.

Inferencing :-



$$SI + d_1(DT_1) + d_2(DT_2) + \dots + d_n(DT_n)$$

\downarrow

[base model]

acc to wt path it is going in a (DT) , we take the avg of those residuals of that particular (DT) split.

\Rightarrow This is also a Black Box model

SVM :- (it is like linear regression)
[Support Vector Machine]

DON'T
FORGET



note: ① if there are no overlaps (or) if there are no points (values) in b/w the hyperplane and Marginal plane.

⇒ then we call it Hard Marginal plane.

② if there are overlaps (or) if there are points in b/w the hyperplane and marginal plane

⇒ then we call it Soft Marginal plane

⇒ we focus on creating this Marginal plane with Maximum Distance, Even though there are some Errors, we consider them by solving it with the help of a Hyperparameter.

How to construct those lines?

$$y = mx + c$$



$$ax + by + c = 0$$

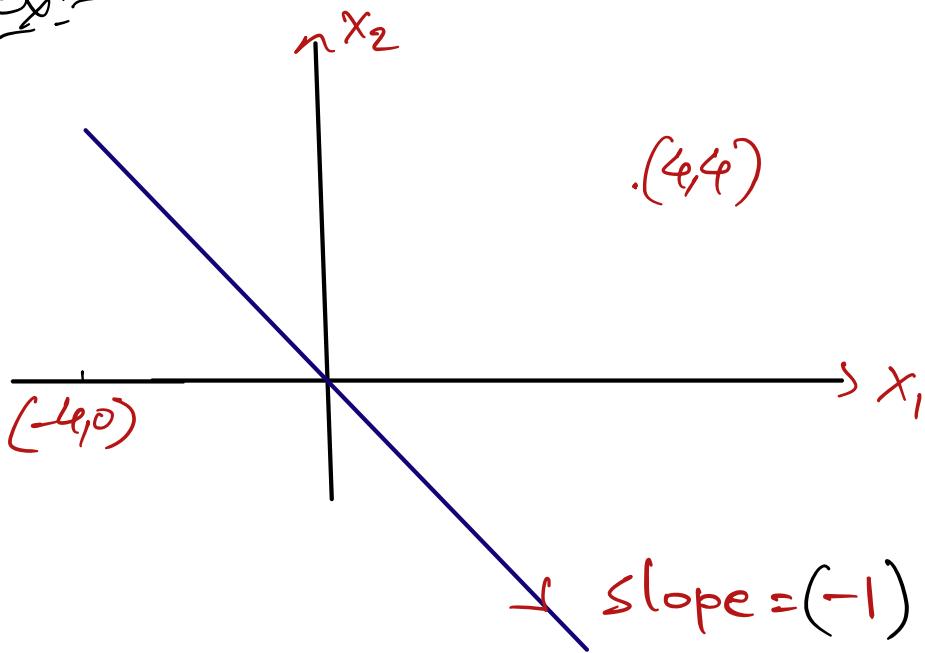
$$y = -\frac{c}{b} - \frac{ax}{b}$$

$$\boxed{m = -a/b}$$

$$y = w_1 x_1 + w_2 x_2 + \dots + b$$

$$\boxed{y = w^T x + b}$$

Ex:-



let, the eqn is passing through origin

$$\Rightarrow b = 0 \text{ and}$$

$$\Rightarrow \boxed{y = [-1] \begin{bmatrix} x \\ 0 \end{bmatrix}}$$

$$\boxed{w = -1}$$

lets, take the $(-4, 0)$ point

$$\Rightarrow y = \begin{bmatrix} -1 \\ 0 \end{bmatrix} (-4 \ 0)$$

$$= 4 \Rightarrow +ve \ Value$$

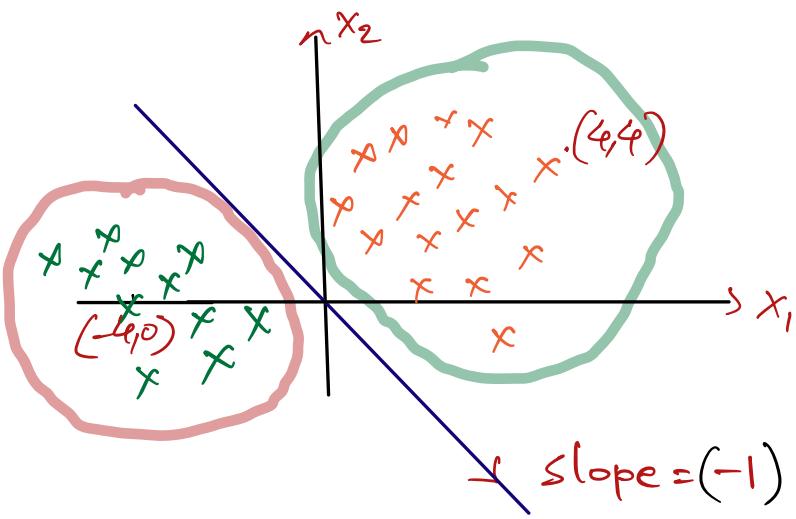
So, all the values, below the line $(0,0)$ in the area of $(-4, 0)$ point are considered to be **+ve Values**

lets, take the $(4, 4)$ point

$$\Rightarrow y = \begin{bmatrix} -1 \\ 0 \end{bmatrix} (4 \ 4)$$

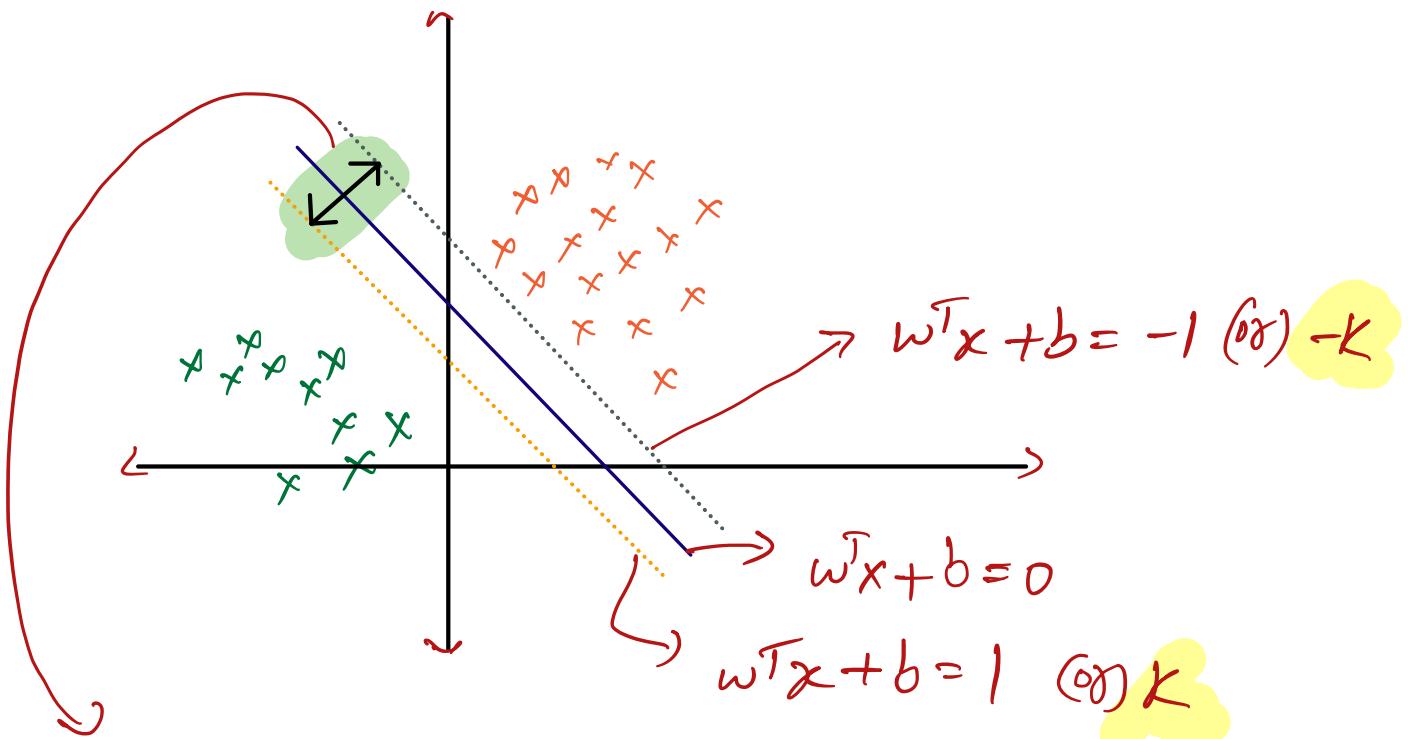
$$= -4 \Rightarrow -ve \ Value$$

So, all the values, above the line $(0,0)$ in the area of $(4, 4)$ point are considered to be **-ve Values**



we divided them,

How to create marginal plane?



our aim is to increase this distance, if we can increase this dist, it means our model is performing well.

$$\begin{aligned} w^T x_1 + b &= 1 \\ w^T x_2 + b &= -1 \\ \hline (-) w^T x_2 - (+) w^T x_1 &= 2 \end{aligned}$$

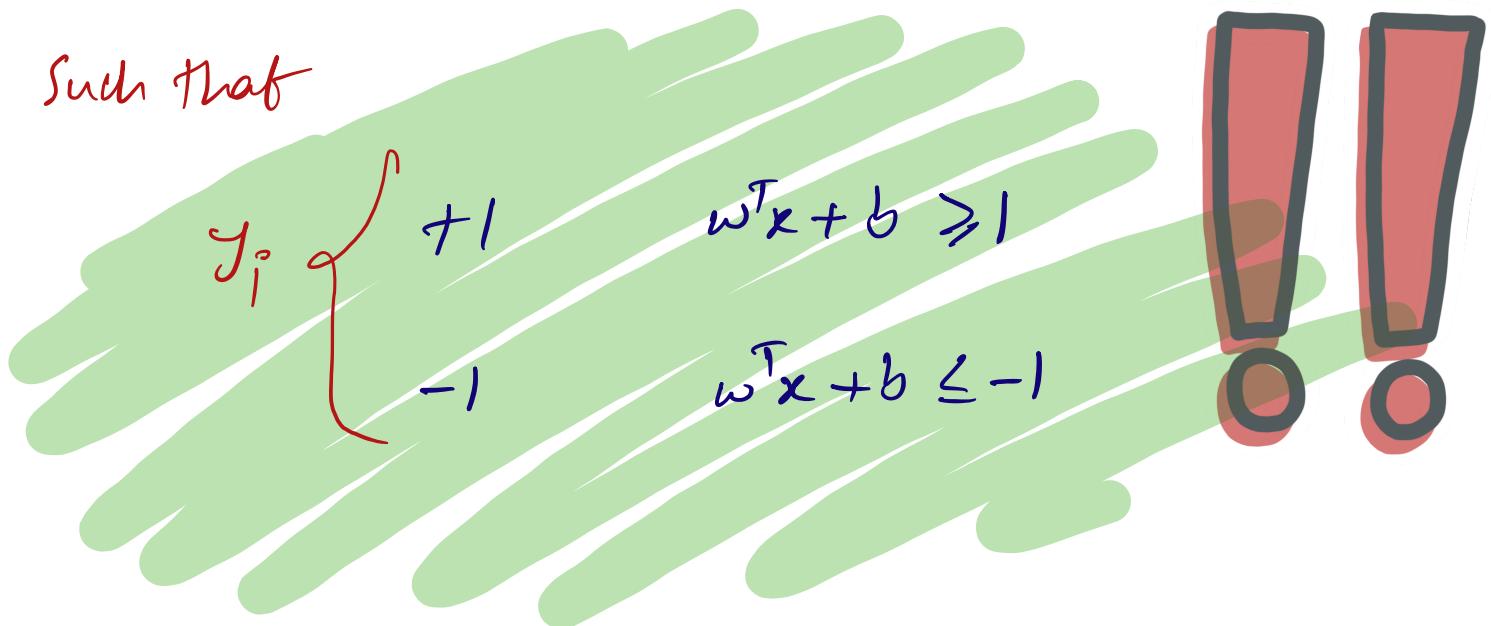
Maximize cost function $\frac{w^T(x_1 - x_2)}{\|w\|}$

To remove the magnitude from a vector, we divide it with its unit vector

Note: our aim is to maximize the value of $2\|w\|$ by updating (w, b) value

→ if we can maximize, then our marginal plane will spread and dist ↑

Such that



Major Aim

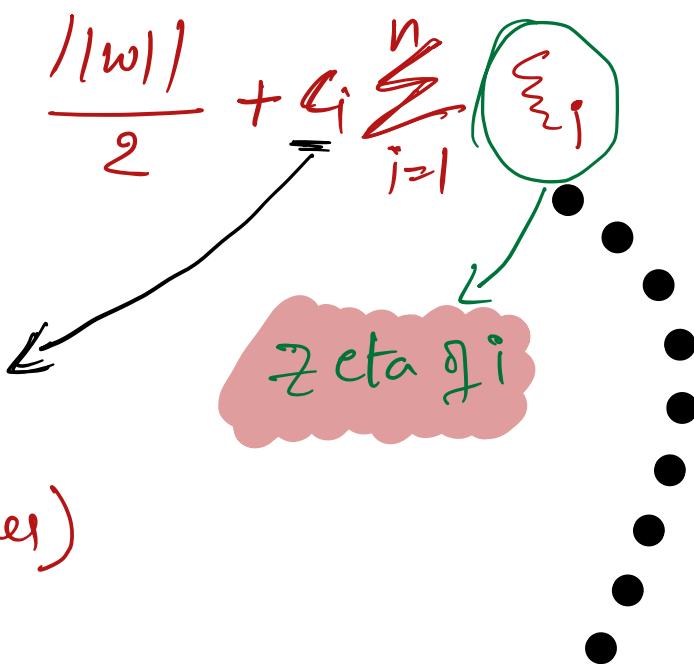
$$y_i * (w^T x_i + b) \geq 1$$

for correct point

for minimize

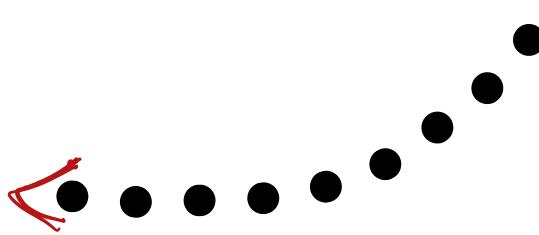
⇒

$$\min_{(w,b)} \frac{\|w\|}{2} + \sum_{i=1}^n \xi_i$$



it says, how many errors we can have .
(like points in opposite class)

Summation of the distance of the wrong data points.



↳ distance from its ext marginal line to the location where the point is.

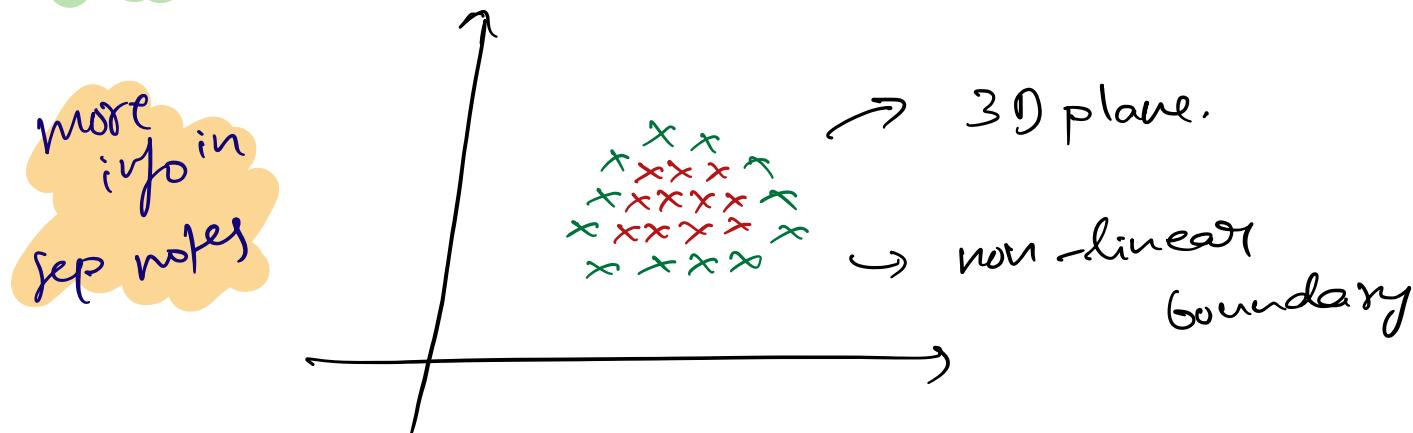
i.e., distance from related marginal line to the steeper class location of that point

$$c_i \sum_{j=1}^n \xi_j$$

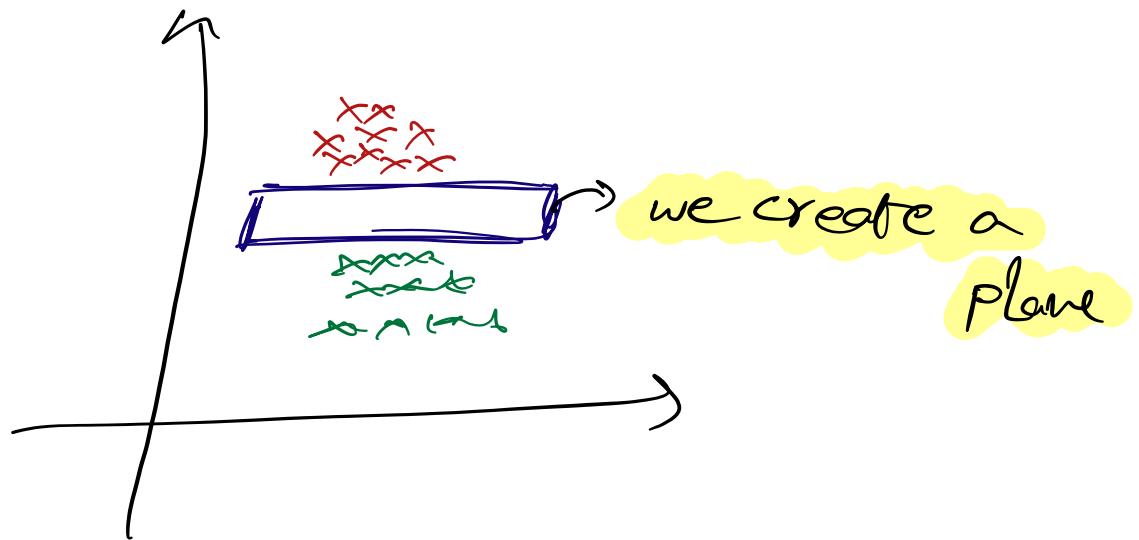
→ this is a hyperparameter used to minimize the errors.

⇒ we can't remove all the error, hence we are ok with few points in wrong location.

SVM Kernel:-



we convert it by pushing the points,
to separate them.



SVR: gives us the flexibility to define how much error is acceptable in our model and will find an appropriate line to fit the data

Minimize:

$$\text{MIN} \frac{1}{2} \|w\|^2 + c \sum_{i=1}^n |\xi_i|$$