

①

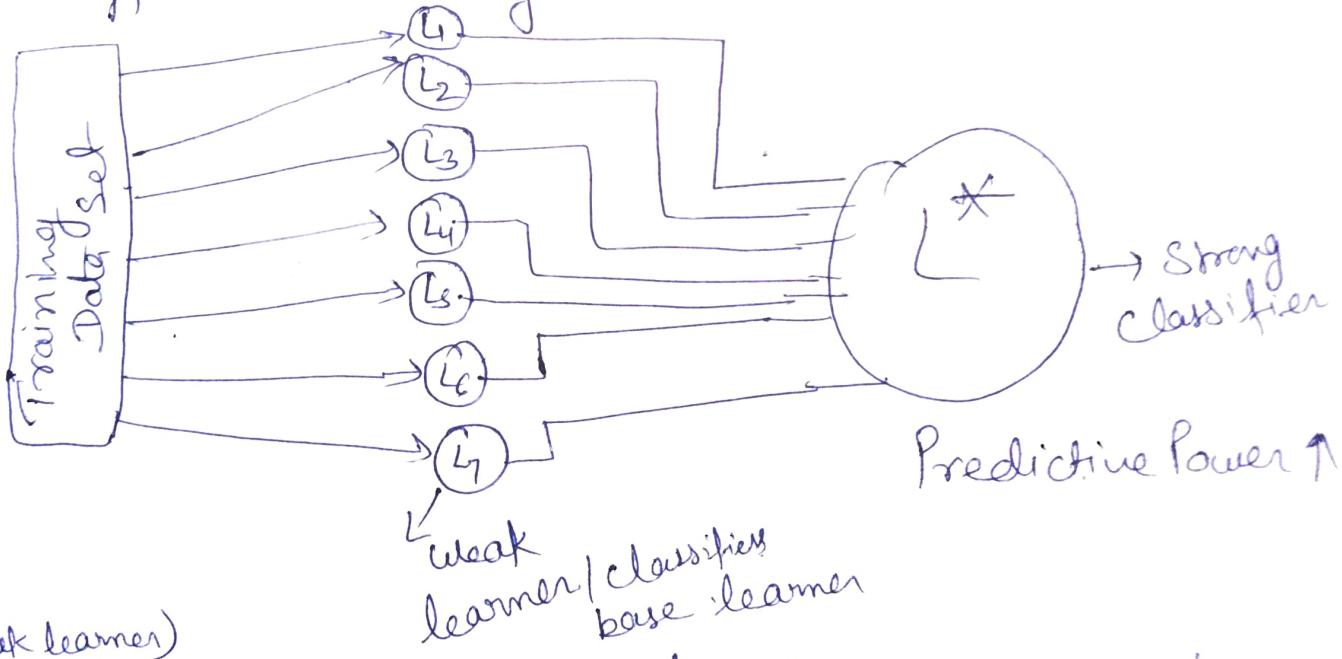
Ensemble learning →

In ensemble learning we will

Combine the output of multiple classifier in order to get better prediction and or classification accuracy.

And under Conditions , Where the classifier output are independent of each other and make errors in an independent manner, it is possible that by Combining the output of Several classifier . we get a resulting classifier which is better than any of the Constituent classifier .

- Different algorithm
- Different Training data set



(weak learner)

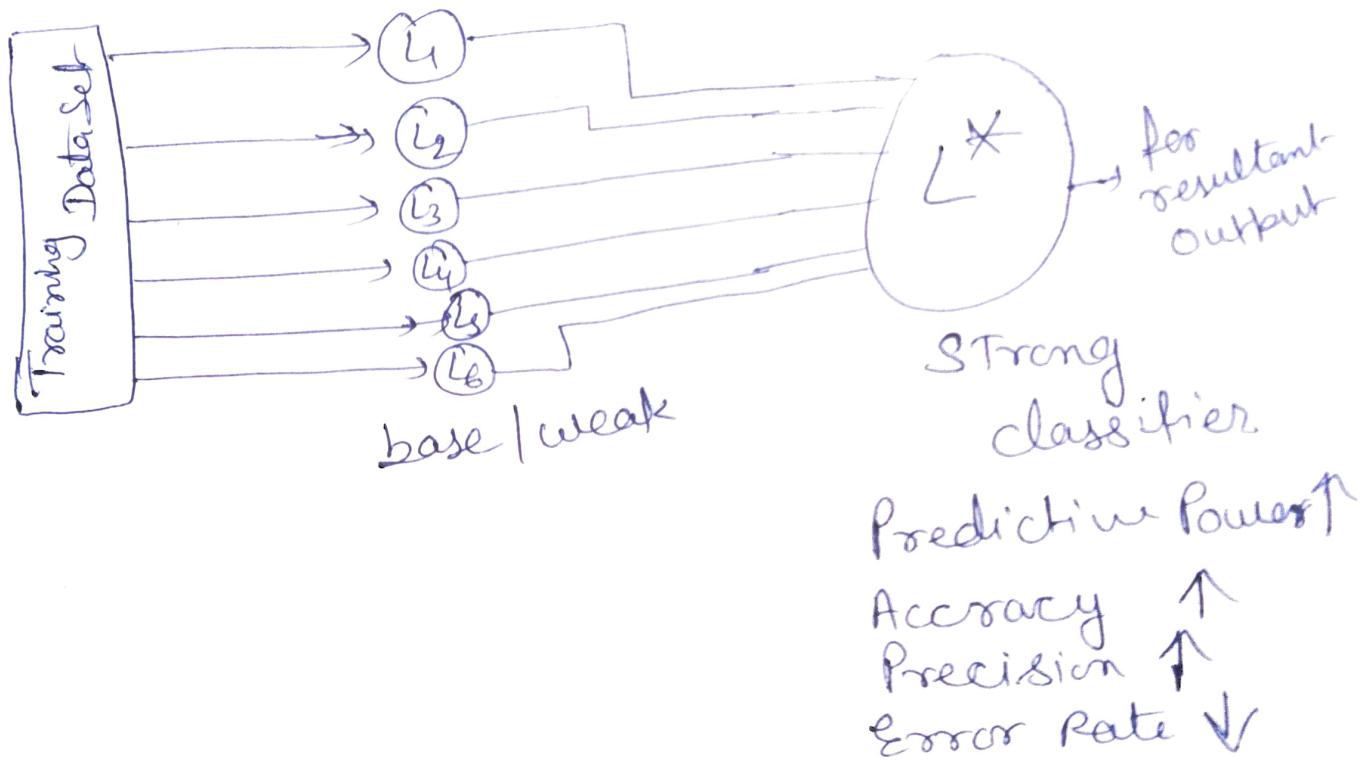
they can use different algorithm to maintain the diversity of their output.

This is Called Heterogeneous Ensemble because all base learners are using different - different algorithm.

②

if all base learners use same algorithm
diversity of output will be low.

To get diversity we can use different-
data Set for each ^{base} learner

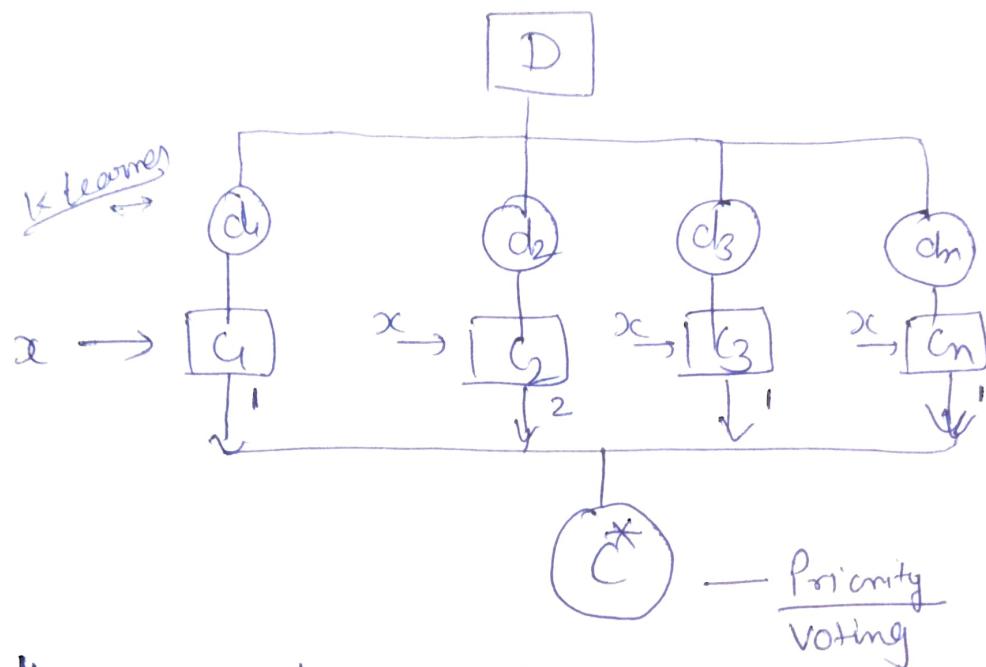


Two Different Ensemble learning Methods

- 1) Bagging
- 2) Boosting.

Bagging → It is also known as BootStrap aggregation.

first of all the samples are generated such that the samples are different from each other we sample D_1 , D_2 , D_3 ..., D_k these are different Data Subsets which is sampled from original Data. They may have some in common



~~but~~ they also will have different values.

To get d_i Draw Random Sample with replacement ..

for Example

If we desire each of these D_i we will have m Examples we randomly draw m samples from D with replacement ~~me~~ that is we can have the same instance repeated several times and some instances may not appear in a particular D_i and so these D_i can be different.

We use bagging technique for these K learners and train a base learner with each and combine the output by voting.

Boosting \rightarrow When we try to get D_i we have to sample from D . But in this we assign a probability with each instance. So every instance is assigned a probability.

Suppose, to start with there are n instances and initially we may have ~~we may give equal probability to all the instances~~ so, each of ~~none of them has~~

In the first iteration, we choose D_1 based on this uniform probability distribution; and ~~we~~ get D_1 by sampling. Now on D_1 we train a Model (M_1) for classification. and this model will apply to all these instances.

Now this Model may have 100% accuracy on these instances. that is the case then it is a very good Model & we no need to really proceed. So it will have when all of them it will do it correctly.

But in Most of time this Model will label some of these instances correctly and some wrongly in D .

	1	D_1
D_2	2	X
X	3	✓
X	4	✓
X	5	X
X	6	X
X	7	✓
	8	✓

D_1 labels.

Some Correctly and Some wrongly.

Now we want the next Model to have 5 independent errors and we want the next Model to do well on those instances where the current Model has not done well.

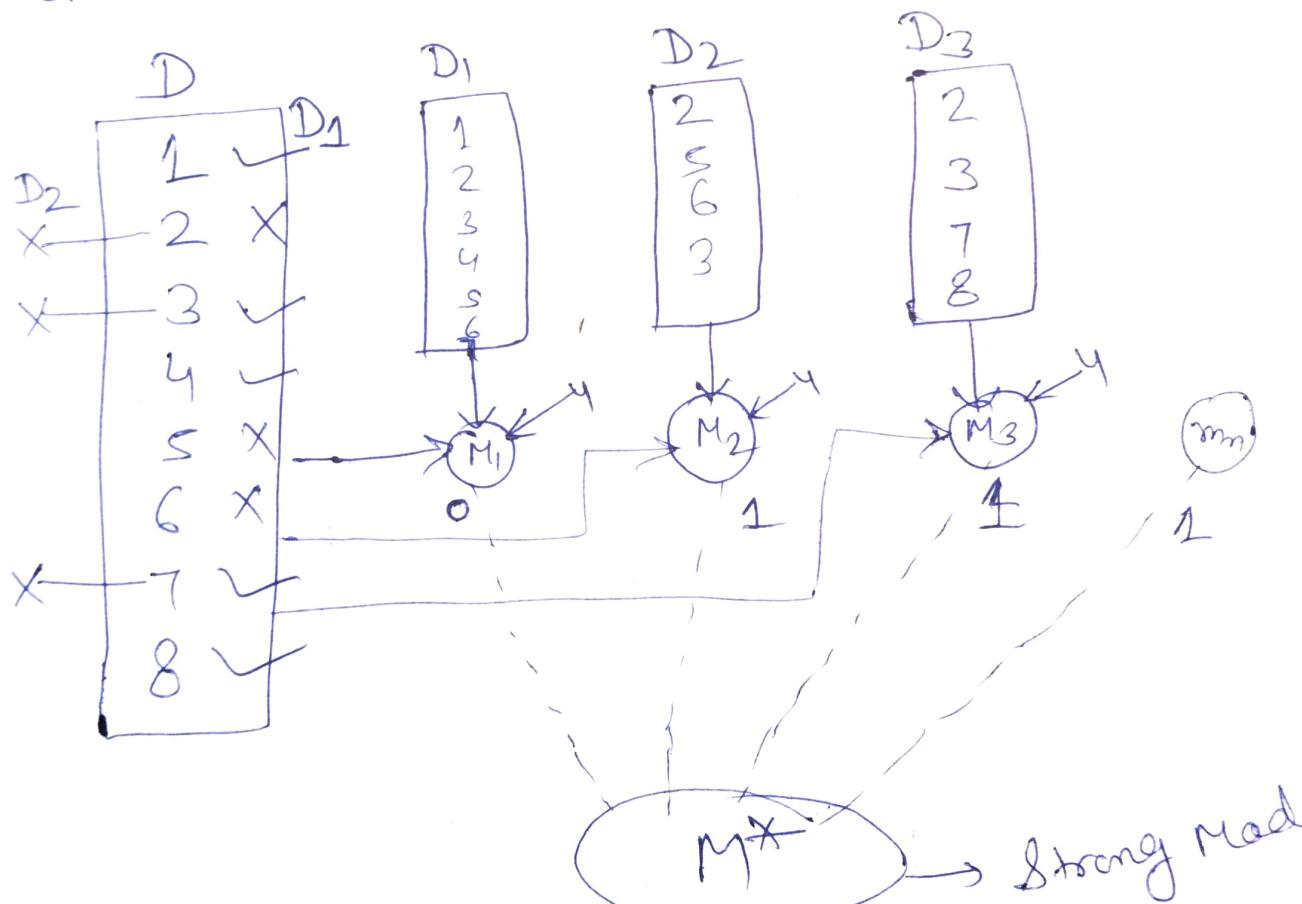
So, now we will change the Probability distribution. And more specifically we identify those instances which were wrongly labeled by Model one, and we will increase the probability of these instances. Those which were correctly labeled we will proportionately reduce the probability.

We want to hide the probability of those instances which we were not correctly classified by Model M_1 . So we will get the same dataset, but now the probability distribution will change. Now we will now sample D_2 according to this probability distribution. and we will get a sub sample and this sub sample is more likely to contain 2, 5, 6 than 1, 3, 4, 7, 8, so we get D_2 and on D_2 we apply Model we apply learning algorithm and get the Model M_2 . Now again we apply M_2 on these training instances and which ones D_2 does error on

Suppose D_2 wrongly classifies 2, 7, 3 and rest it does correctly so, so we will increase the probability of 2, 3, 7

beyond that was the current probability and reduce the probability of the rest to keep the sum of the Probability equal to 1. ⑥

So based on that we will find D_3 on which we will train Model M_3 like this we will go on, we will get the different Models a different data and the different Models and then we will combine these different outputs of different Models by voting and Kth.



class → 0

→ 1 Voting ↑
1

$$Y=1$$

Random Forest Algorithm

7

- * Random Sampling of Training data when building the trees - this first Randomness.
- II Randomness → Random A Subset of input features when splitting at Nodes.
for example

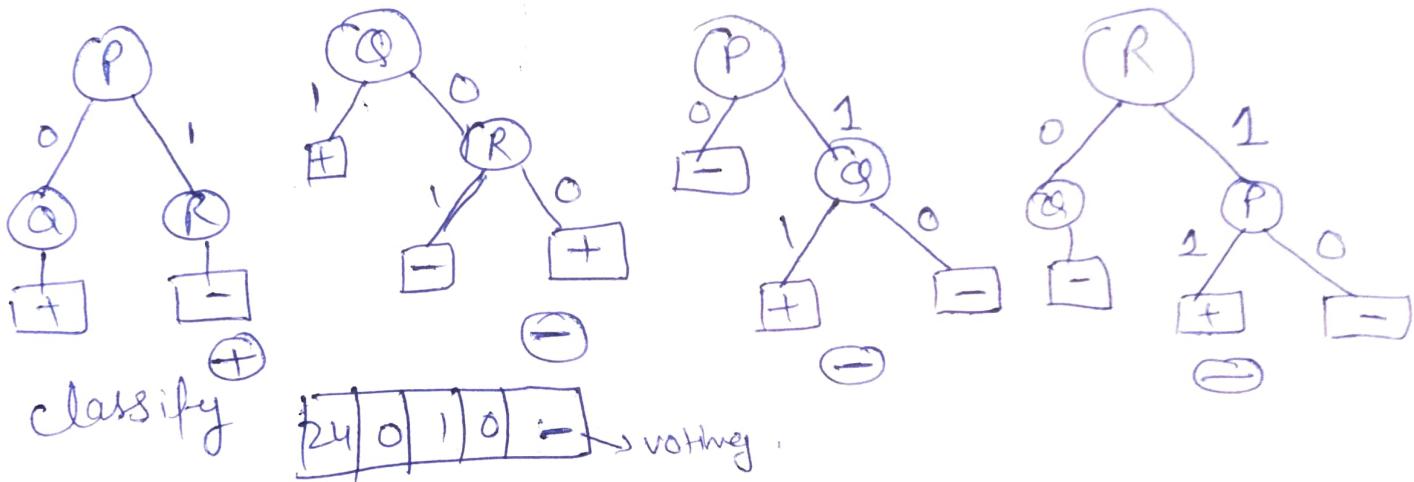
Dataset -

ID	P	Q	R	Class
20	0	0	1	+
21	1	0	0	-
22	2	0	0	-
23	1	1	2	+



Bootstrap Sample

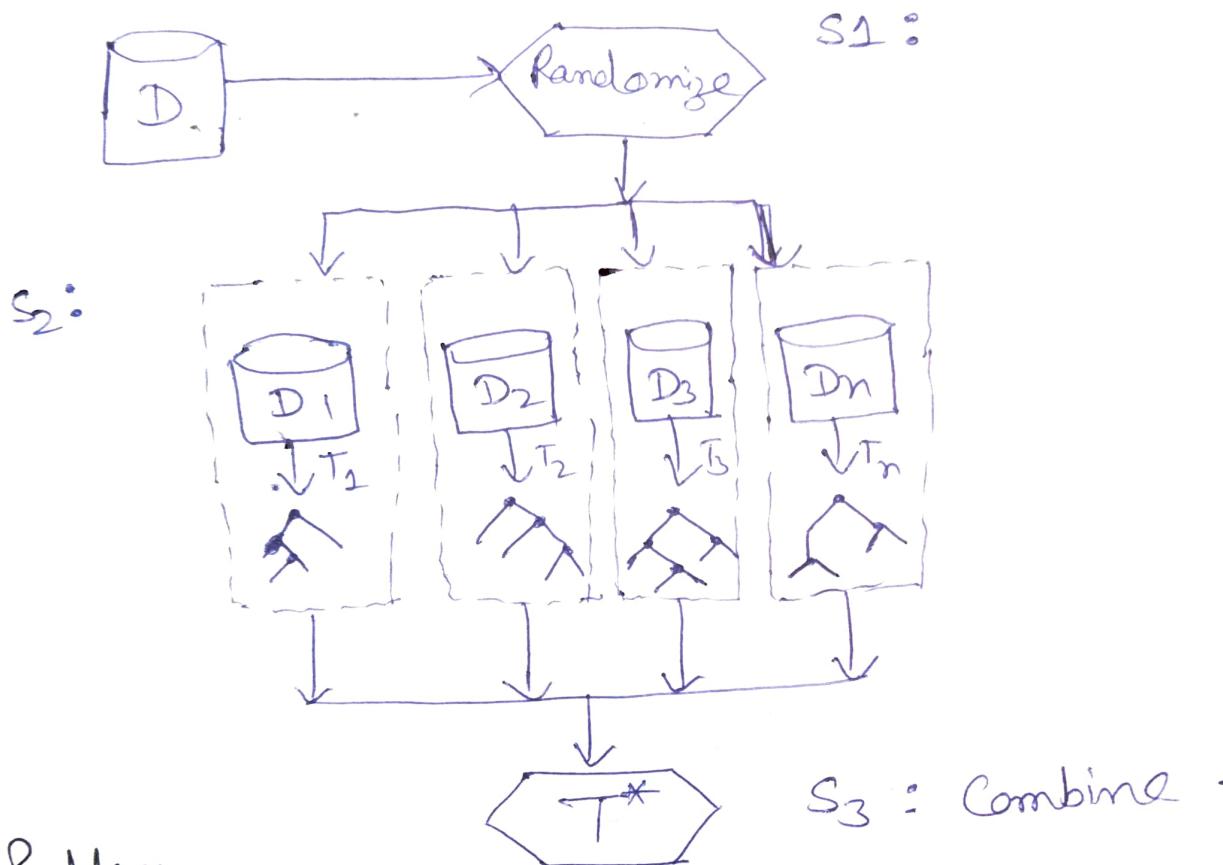
ID	P	Q	R	Class
20	0	0	1	+
22	2	0	0	-
20	0	0	1	+
22	2	0	0	-



for creating a bootstrap sample we can have duplicate entries or duplicate observations - so this we have random sampling of the data means we have taken randomly. In some random fashion. Second the random subset of input feature selection when splitting at the nodes now. P, Q, R are the instances or attributes which we have for this data set or bootstrap sample. so this we create one instance which have highest Information Gain as you mentioned.

And rest two will become the root nodes are 8 intermediate nodes.

So we can create n number of bootstrap sample or n number of Decision tree and prediction are combined using this voting mechanism.



Problem

Day	Outlook	Humidity	Wind	Play
D ₁	Sunny	High	Weak	No
D ₂	Rainy	High	Strong	No
D ₃	Overcast	High	Weak	Yes
D ₄	Rain	High	Weak	Yes
D ₅	Rain	Normal	Weak	Yes
D ₆	Rain	Normal	Strong	No
D ₇	Overcast	Normal	Strong	Yes
D ₈	Sunny	High	Weak	No
D ₉	Sunny	Normal	Weak	Yes
D ₁₀	Rain	Normal	Weak	Yes
D ₁₁	Sunny	Normal	Strong	Yes
D ₁₂	Overcast	High	Strong	Yes
D ₁₃	Overcast	Normal	Weak	Yes
D ₁₄	Rain	High	Strong	No

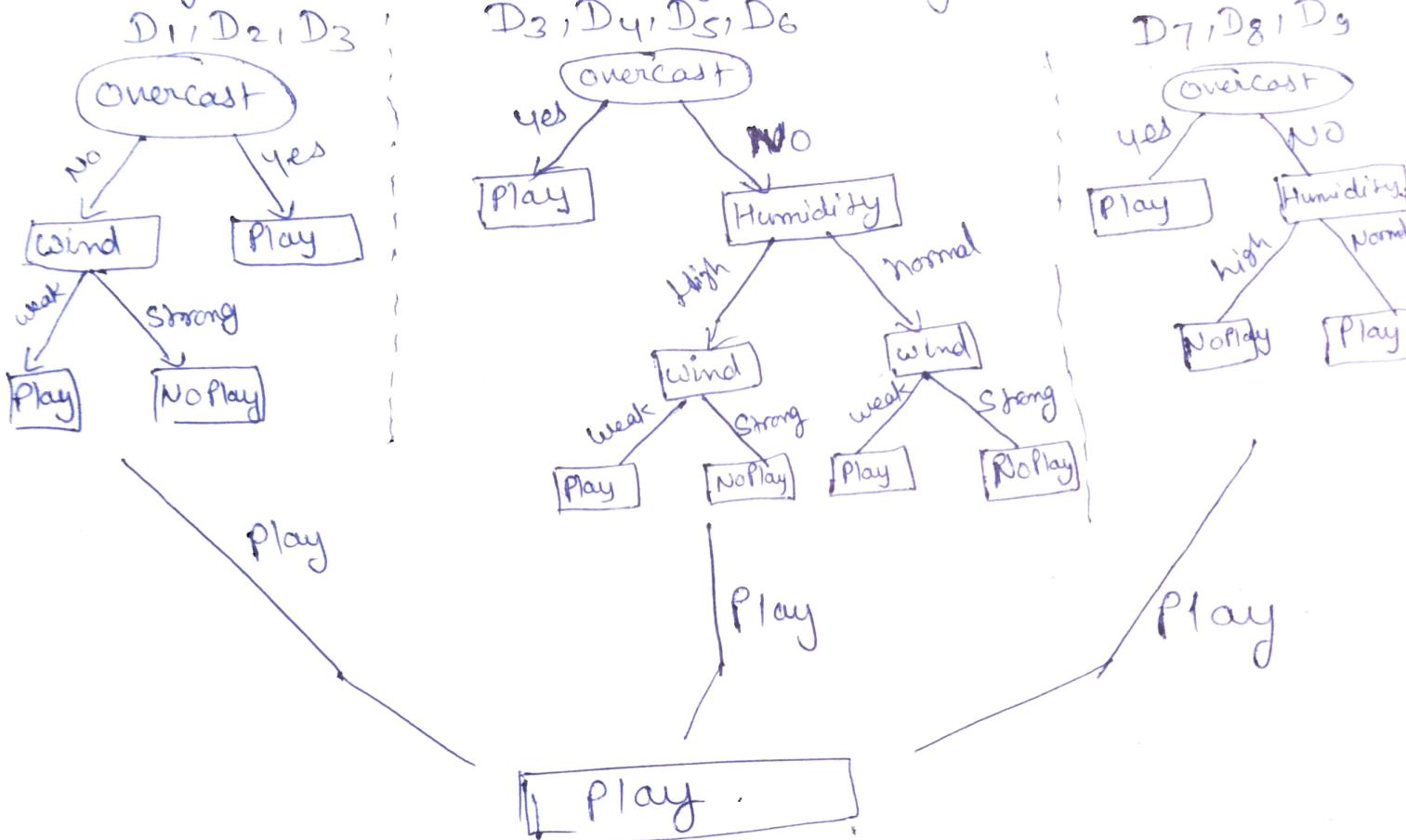
Outlook - Rain
Humidity = High
Wind - weak
Play = ?

Now using random forest - we need to predict whether the game will happen if the weather condition is above mentioned.

* We divide the data into smaller subset.

(9)

* Every subset ~~need~~ may or may not be overlapped.



features → Most accurate learning algorithm

- Runs efficiently on large database
- Perform implicit feature selection
(automatically) collect some random feature
- Can be easily grown in parallel so it
- ~~to~~ reduce computation time

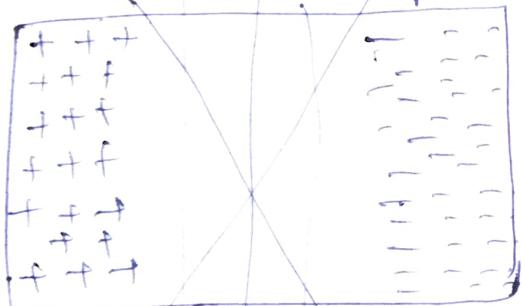
Support Vector Machine (SVM) →

Applications:

- ① Hand written character / digit recognition
- ② it Perform well on high Dimensional Data and Solve the Problem of High Dimensionality.

Support vector:

↓
it means Examples or training data
Subset of training data and used to represent decision boundary.



→ These are hyperplane.

These can be infinite in number

* we can choose most optimal one.

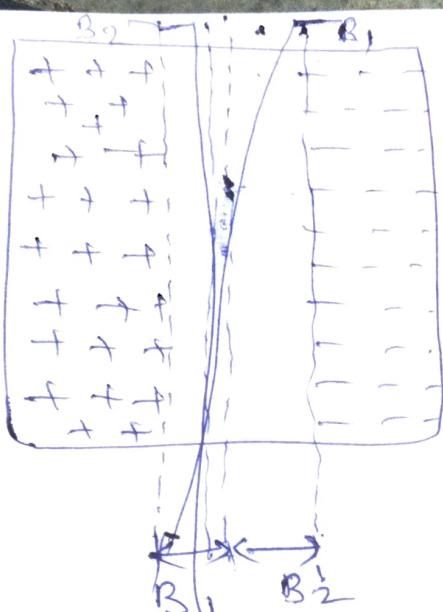
* Training Error on this 0

The chosen hyperplane may not perform in future. It's not equally performed well on any unseen Examples.

Suppose we choose one hyperplane and we use it another day for classification then it can not equally separate all the unseen future instances. So this is the drawback of this.

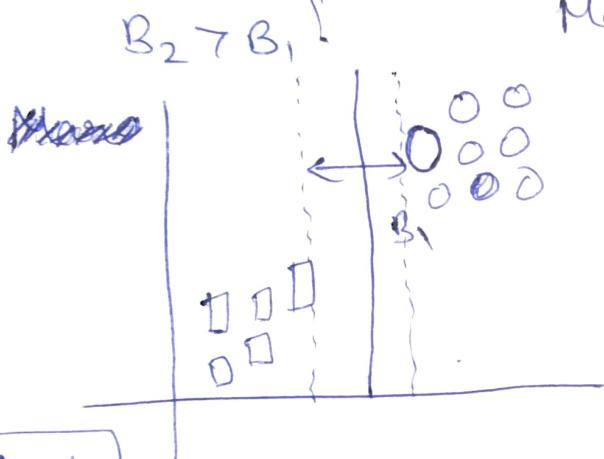
for choosing the hyperplane we use maximal margin hyperplane.

(11)



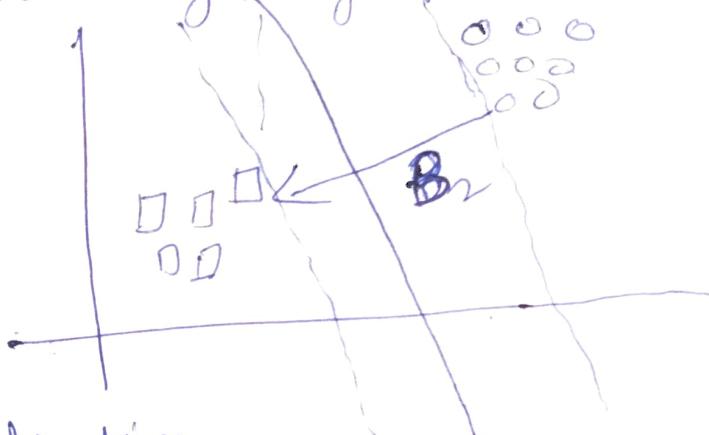
Margin $\uparrow \rightarrow$ better Control on generalization error

Margin $\downarrow \rightarrow$ Model will fail.
"victim of overfitting"



<input type="checkbox"/> yes
<input type="radio"/> no

Maximal margin hyperplane



Linear classification.

Maximal Margin gives better future,
Prediction.

Margin \uparrow

Misclassification \downarrow
Accuracy \uparrow
error \downarrow

Accuracy yes

$\leq 50\%$.

binary class \rightarrow line

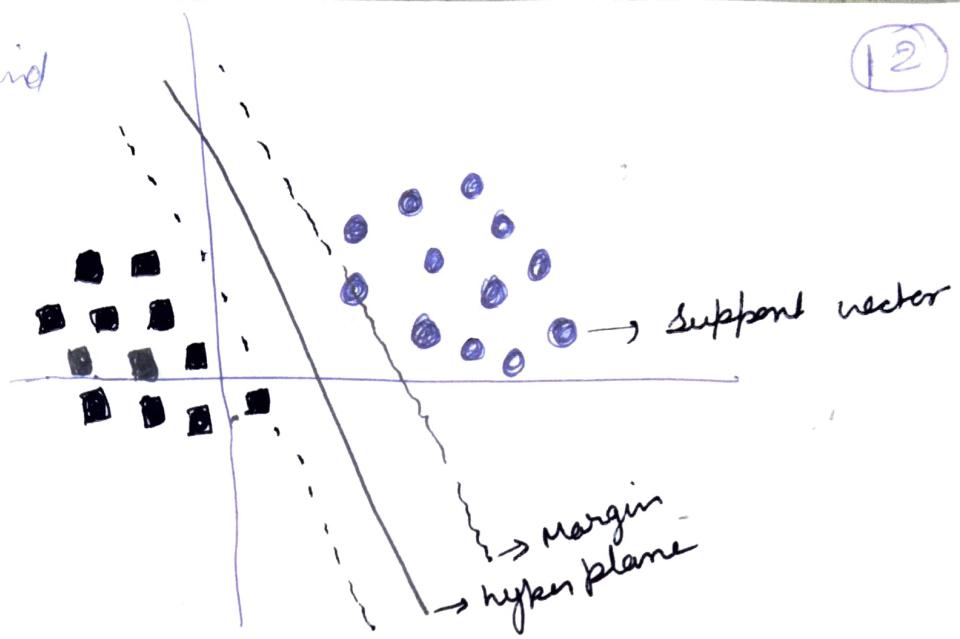
3D \rightarrow plane

4D \rightarrow hyperplane

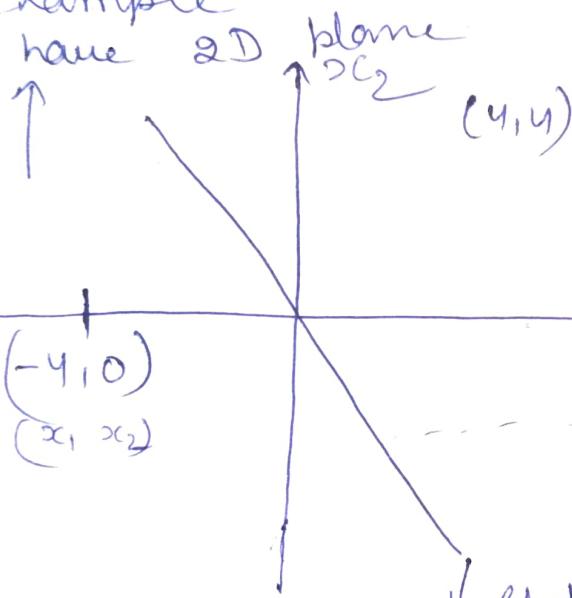
SVM Kernel is used, to convert dimensions
of plane. So it convert 2D to 3D, 4D
increase

(12)

We have to find
out Plane



for Example
we have 2D



Equation of line

$$y = mx + c$$

$$m = -1$$

$$c = 0$$

slope = -1

$$y = mx + c$$

if we convert this line equation in to
vector. then

$$y = w^T x + c$$

and if it is hyperplane then

$$w^T x + b = 0$$

i.e. if any point exists in the hyperplane
then value of sum is zero.

Here
 $x = (x_1 \ x_2)$
 $w^T = w$ transpose of w

$|w|$ = Vector

Now

$$Y = \mathbf{w}^T \mathbf{x} + C$$

$$m = -1$$

$$C = 0$$

(13)

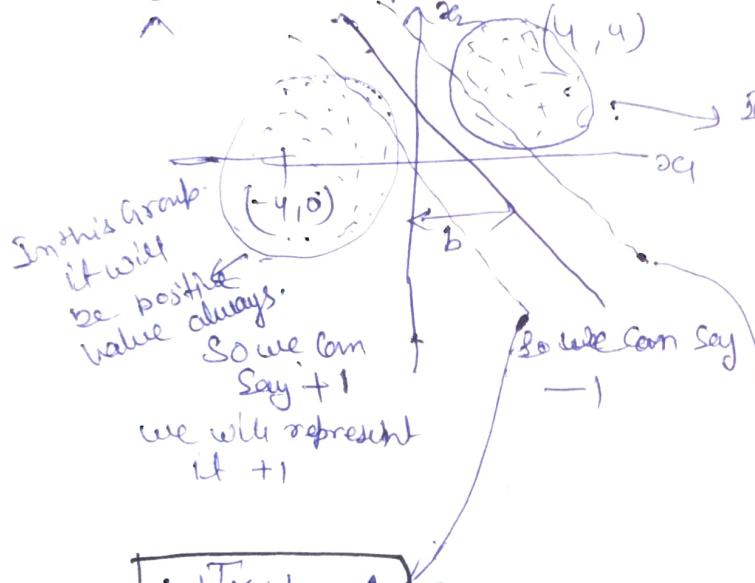
$$Y = [-1] [-4, 0] + 0$$

$$\mathbf{x} = [x_1, x_2]$$

$$\mathbf{w}^T \begin{bmatrix} -1 \\ 0 \end{bmatrix} \begin{bmatrix} m \\ C \end{bmatrix}$$

$$\underline{Y = Y}$$

if we will put any value in this equation so it will be positive in the left side.



In this graph
it will be negative always.

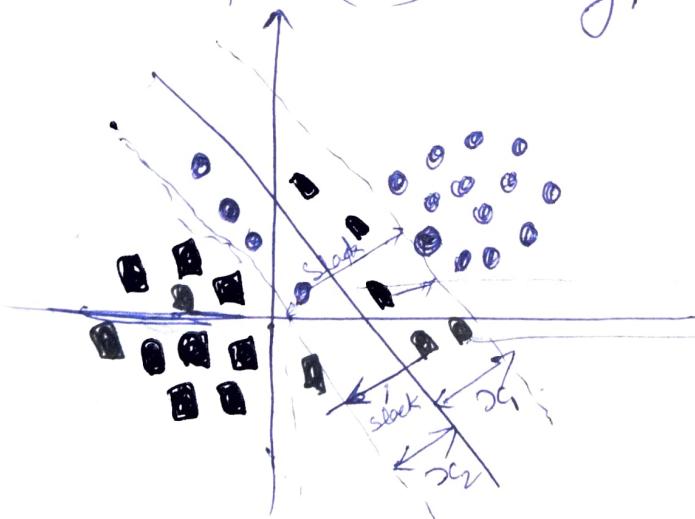
$$Y = \mathbf{w}^T \mathbf{x} + C$$

$$= [-1] [4, 4] + 0$$

$$\underline{Y = -4}$$

$$w^T x + b = -1$$

b is distance from $(0,0)$ to hyperplane



we have to find out $\|x_2 - x_1\|$
because points are in opposite direction

$$\cancel{w^T x_1 + b = -1}$$

$$\cancel{-w^T x_2 + b = +1}$$

$$\frac{\cancel{-w^T x_2 + b = +1}}{\cancel{w^T (-x_2 + x_1)} = -2} \Rightarrow w^T (x_2 - x_1) = +2$$

$$\Rightarrow w^T (x_2 - x_1) = 2$$

In this situation
misclassification
of these
point so
Model will be
overfit.
we have to select
optimized hyperplane

$$w^T(x_2 - x_4) = 2$$

we have to find out $x_2 - x_4$, so we have to replace w^T for that

$$\frac{w^T}{\|w\|} (x_2 - x_4) = \frac{2}{\|w\|}$$

$$\frac{w^T}{\|w\|} = 1$$

Nom of w

$$(x_2 - x_4) = \frac{2}{\|w\|} \quad \uparrow \uparrow \uparrow$$

for all $w \neq 0$

$$(w^*, b^*) \text{ max } \frac{2}{\|w\|}$$

such that

$$\begin{cases} y_i = +1 & w^T x + b \geq 1 \\ y_i = -1 & w^T x + b \leq 1 \end{cases}$$

→ this is optimized function.

$$y_i * w^T x + b \geq 1$$

↓ it is always true if it is not true then misclassification.

and if we have to find min of any maximum value then we have to do reciprocal of max.

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

hyperplane
slack variable

No of point → How many error you want in misclassification to allow.

final equation.

this will be regularization.

to avoid overfit.

Slack variable =
value of error

it is a distance between point to margin.