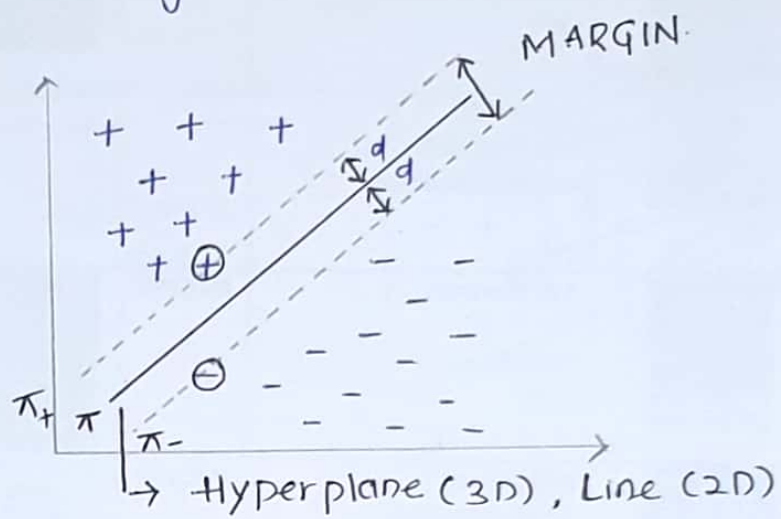


## Support Vector Machine Algorithm :-

It is also called mostly SVM. In this our task is to find a line that best separates the positive points from negative points as widely as possible (Simply make a maximum margin line), it is a set of supervised learning methods used for classification, regression and outlier detection, the advantages of SVM's are - Effective in high dimensional spaces.



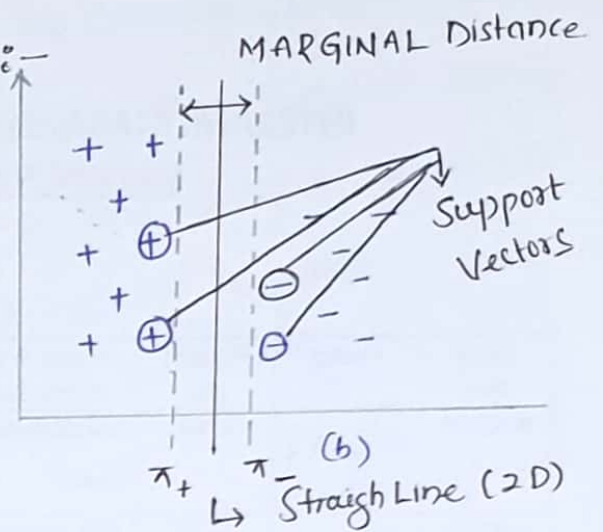
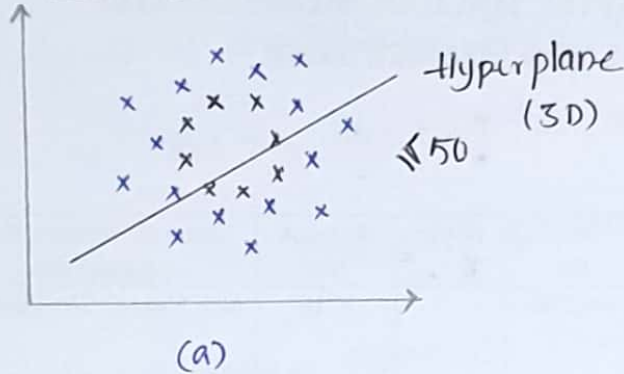
From the figure we can say that a line (2D) which separates two data points into two classes as +ve's, -ve's, and also makes two marginal lines ~~also~~ and this widely separated parallel lines having some distances, these lines are passing through one of the nearest points maybe positive or negative, the difference between logistic regression and SVM is both have separated the points but SVM has some widely difference distances.

The distance between ( $\pi_+$ ) line to hyperplane is " $d_+$ " and distance between ( $\pi_-$ ) line to hyperplane is " $d_-$ ".

Summation of this two distances is called as "Margin" (2)

The role of ~~margin~~ margin is to divided a class in better way to ~~Simple~~ tell us the class of a data points - for suppose if a new data point comes and lying in between hyperplane and one of margin then we simply understood the point belonging to that class.

Linear Separable & Non Linear Separable :-



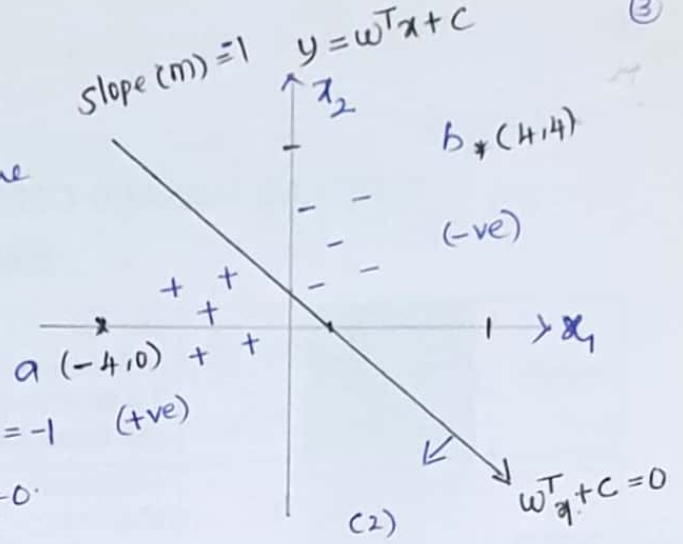
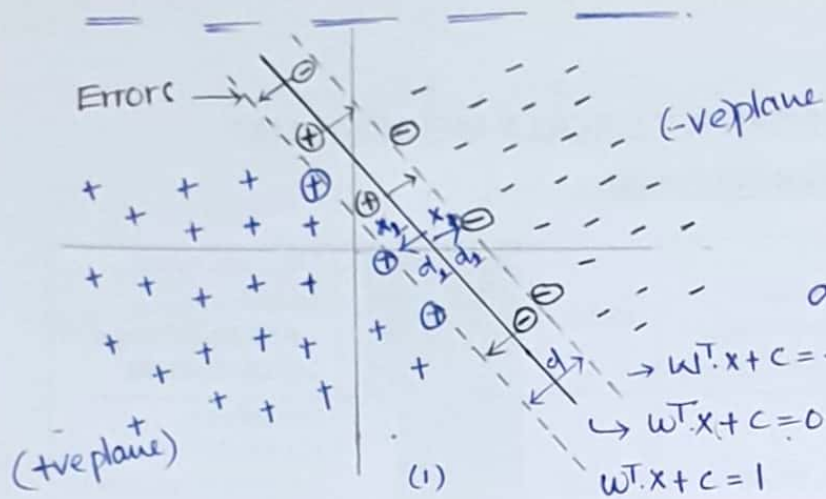
If we observe fig (a) all the points are overlapping with each other there we can't separate data points by making a plane or hyperplane. So we can say this is a non-linear separable planes, if solve like this problems we have to use kernel distribution.

from figure (b) we easily separate the data points by making a straight line this is called as linear separable line.

After classified data points the nearest data points from any class, passes through marginal lines are called as Support vectors as shown in figure if new data point come that is also.



# SVM MATHS INTUITION :-



the main diff between logistic regression and svm classification is marginal lines this marginal lines how we created we discuss here from fig(2) if we observe a line we are making that separates having slope = -1 and the line eqn is  $w^T x + c = 0$  and we know  $y = w^T x + c$ . let consider two points (4, 4) and (-4, 0) as shown in graph.

for Suppose point (a) i want 'y' value -

So here  $m = -1$  ; Constant (c) = 0 because passes through origin.

$$\text{Then } - y = \underline{w^T x} + 0 \Rightarrow$$

$$y = \begin{bmatrix} -1 \\ 0 \end{bmatrix} \begin{bmatrix} -4 & 0 \end{bmatrix}$$

$$w^T \quad x = 'a'$$

$\therefore$  Dot matrix

$$y = 4, \therefore \text{When ever we do Dot matrix we get Scalar value only.}$$

$\therefore$  y value is +ve

Note : When ever we get a line the below mentioned points all are always to be positive and Similarly the opposite Co-ordinate values always to be negative.

$$\begin{bmatrix} -1 \\ 0 \end{bmatrix} \begin{bmatrix} 4 & 4 \end{bmatrix}$$

$$y = -8,$$

after find out  $y$  value we get two groups that are (+ve), (4)  
 (-ve) classes, if I consider line below group as "+1" because  
 it gives positive points value may be differe. and above as "-1"  
 for easy calculation.

from figure (c) let consider nearest points to the line (ie  
 shortest distance from line to point on below side, above side of  
 a line. and make a dot point of line parallel to 'y' line that indicates  
 marginal lines. and distance b/w them is marginal distance.  
 if we observe clearly it, create the positive plane & Neg. plane

So from that the line eqn's are -  $y = w^T \cdot x + c = 0$ .

below Hyperplane is -  $w^T \cdot x + c = 1$

above Hyper plane is -  $w^T \cdot x + c = -1$

In graph I mentioned points as  $(x_1, x_2)$  which are nearest points  
 in different planes of different marginal lines. and the distances are  
 $(d_1, d_2)$  respectively. -

$$w^T x_1 + c = -1 \quad \text{--- (1)}$$

$$w^T x_2 + c = +1 \quad \text{--- (2)}$$

$$\frac{-}{-} \frac{-}{-} \frac{-}{-}$$

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

from this two eqn we can get total distance b/w two parallel  
 planes. to plane, now I get the distance but I want to remove  
 the  $w^T$ . In Order to do that we have to use  $\|w\|$  length of  
 vector because we know direction.



Now we have to divide Eqn by  $\|w\|$  both side. (5)

$$\frac{w^T (x_2 - x_1)}{\|w\|} = \frac{2}{\|w\|}$$

→ This is our Optimization Eqn, and we have to Increase this value. (Maximize it)

↳  $w^T$  Magnitude goes one only direction there

Simply  $\Rightarrow (w^*, c^*) = \text{Max} \frac{2}{\|w\|}$  → Optimization function.

①  $\Rightarrow y_i * w^T x_i + c_i \geq 0$

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

HARD MARGIN SVM

from this weq when ever i get  $w^T x + b$  value above 1 always i have to consider it as +ve value similarly below -1 always consider as negative value.

from that eqn we can say the value is always greater than "1" if it is not, then we understand there is a misclassification.

from the "SVM" Eqn we have to minimize the distances for better optimization. In order to do it we have to do reciprocal of it.

$$\text{Max } f = \text{Min } \frac{1}{f} \Rightarrow \text{for minimum value}$$

The Eqn will be  $w^*, c^* = \text{min} \frac{1}{2/\|w\|}$

$$w^*, c^* = \text{Max} \frac{2}{\|w\|} \Rightarrow w^*, c^* = \text{min} \frac{\|w\|}{2} \Rightarrow \text{Hard Margin SVM} \quad \text{--- (2)}$$

for better Optimization we have to include c value and Zeta value where 'c' → Indicates how many errors will be model Included,  $\xi$  → value of the error.

Here final eqn will be -

(6)

$$W^*, C^* = \text{Min} \left\{ \frac{\|W\|^2}{2} + C_2 \sum_{i=1}^n \xi_i \right\} \quad - (2)$$

Where 'c' is hyper parameter and it represent errors, which we are considered i.e. for Suppose if the test data values above the marginal lines that should be consider as error, here we have to do treatment on it, and the errors will be low in value counts, so we can easily low the values of error i.e. negligible by hyper parameter tuning 'c', so initially in Eq (2) model has issue of overfitting we applied hyper parameter as shown in Eq (3) that tune the model, in order to get 'best model'.

So here Eq (1) represent Hard Margin S.V.M formulation

Eq (2) represent Simple S.V.M formulation.

Eq (3) represent Simple SVM formulation with Regularization term where C is hyper parameter.

from Hard Marginal SVM formulation we will get dual form of S.V.M

Dual form of S.V.M :- Above we discussed formulation was the Simple form of S.V.M the alternate method is dual form of S.V.M. which uses

"Lagrange's Multiplier" to solve the Constraints Optimization problem

Note :- If  $y_i * (w^T x_i + c_i) > 0$  then  $x_i$  is a Support vector and when

$y_i * (w^T x_i + c_i) = 0$  then  $x_i$  is not a Support vector.

Eqn will be -

$$\text{Max} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \cdot \alpha_j \cdot y_i \cdot y_j \cdot x_i^T \cdot x_j$$
$$\text{st } 0 \leq \alpha_i \text{ \& } \sum_{i=1}^n \alpha_i y_i = 0.$$

- (4)



from this dual form eq (4) while applying Gradient descent (7)

Eq (4) - represent dual form of Hard Margin S.V.M formulation

if we observe eq 4  $\rightarrow \{x_i^T \cdot x_j\}$  represent dot product of familiar object - this notation is called as "kernel" trick.

kernel trick eqn - 
$$\text{Max} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \cdot \alpha_j \cdot y_i \cdot y_j \cdot k(x_i, x_j) \right\} \quad - (5)$$

where -  $k(x_i, x_j) = x_i^T \cdot x_j$

So from this we can say -  $k(x_i, x_j) = x_i^T \cdot x_j = \{x_i \cdot x_j\}$   $\rightarrow$  This term is

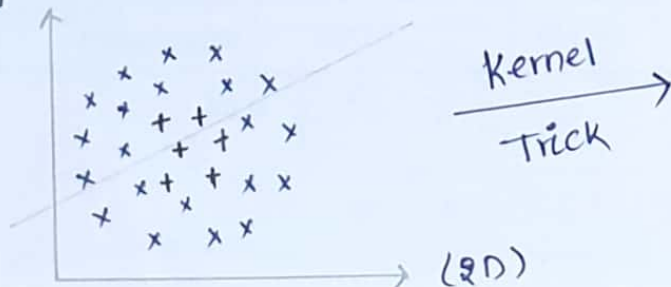
represent linear S.V.M, it is nothing but also called as logistic regr.

because in linear S.V.M the linear line just classify the data Only not mentioned marginal lines. ie simply acted as logistic regression.

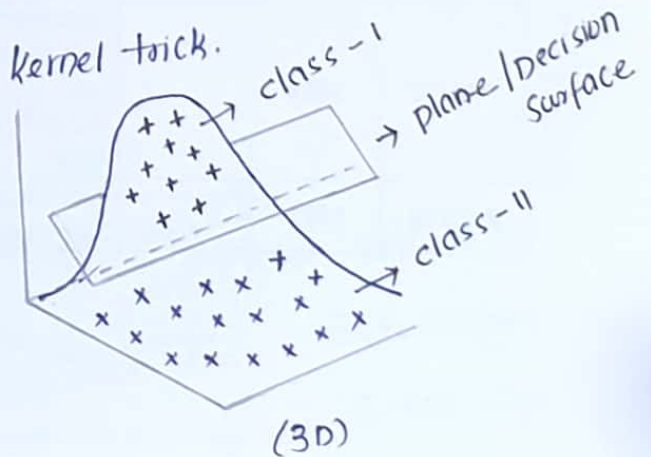
So Instead of using Simple linear SVM most use prefer polynomial or Quadratic kernel SVM which looks like -

$$k = (x_i \cdot x_j) = (x_i \cdot x_j)^2 \quad - (6)$$

So here we just Squaring the eqn of linear S.V.M or we can say we apply degree of '2' in polynomial eqn (or) Quadratic kernel S.V.M, now we plot the graph in (2-D. and 3-Dimension) where we understand how the plane cuts or Separated data points. by kernel trick.



Kernel Trick  $\rightarrow$



By above figure Shows kernel trick applying data and

(5)

eq(6) known as polynomial kernel SVM of degree = 2, we can change degree also in the eqn.

In Generally we use "R.B.F kernel SVM" most of the times, this is an important eqn when ever we don't know to use which algorithm of SVM use in that case we go with "R.B.F kernel SVM".

$$k(x_i, x_j) = \exp \left\{ \frac{-\|x_i - x_j\|^2}{2\sigma^2} \right\}$$

This term 'R.B.F' stands for Radial Basis kernel.

Upto now we discuss SVM eqn in diff-formats such as linear, Simple SVM eqn, Quadratic eqn, polynomial equation with different degrees and finally we using hard margin eqn with kernel tricks to best Optimization. Now we discuss Soft margin eqn, by applying "Lagrange's Multiplier" as similar to Dual form hard margin eqn, the difference is in a hard-margin SVM a single outlier can determine the boundary, which makes the classifier overly sensitive to noise in the data, the result is that Soft margin SVM could choose decision boundary that has non-zero training error even if data set is linearly separable and is less likely to overfit.

$$\begin{aligned} \text{Max } & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \cdot \alpha_j \cdot y_i \cdot y_j \cdot (x_i^T \cdot x_j) \\ \text{st } & 0 \leq \alpha \leq C \quad \& \quad \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

⇒ called as  
Soft Margin  
S.V.M. Eqn.

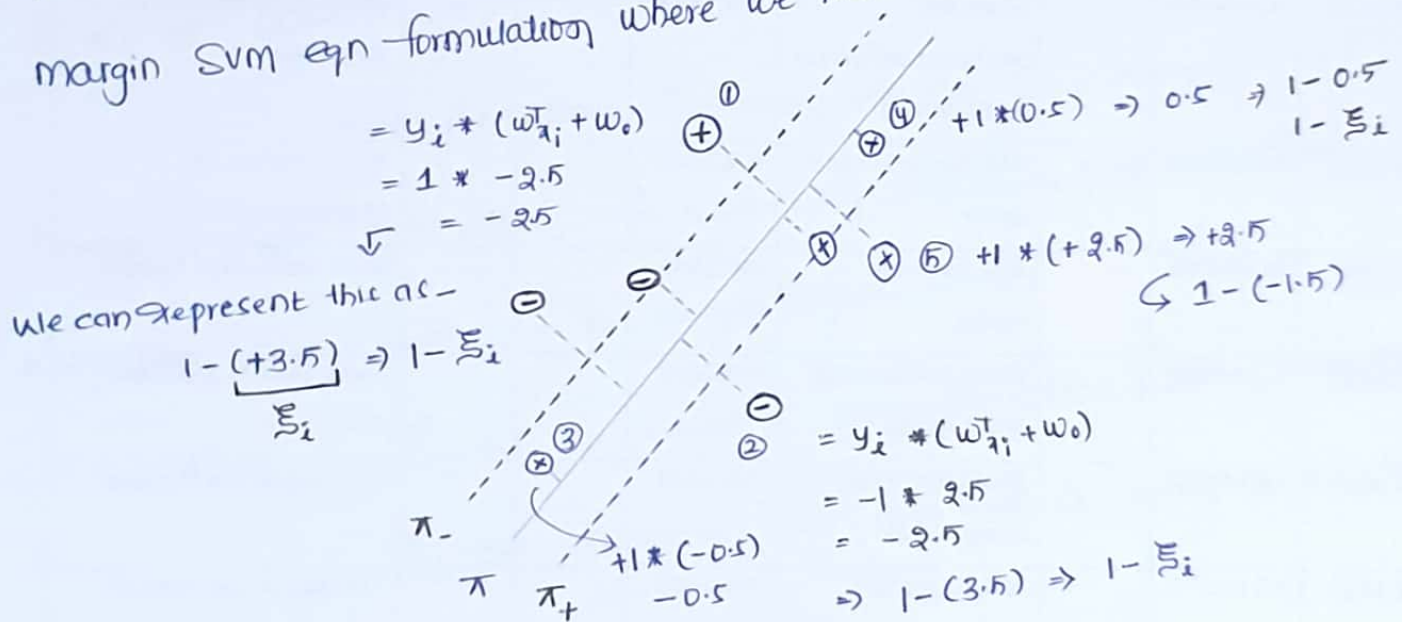


As similar in hard margin eqn what ever changes we done (like kernel trick, Quadratic eqn, polynomial eqn with different degrees) we do hear also, we can replace that term also.

In hard Margin s.v.m we are trying to maximize the margin (Signed Distances) to all the points. ie  $(y_i * (w^T x_i + w_0)) \geq 0 \forall i$  when there is no noisy in data. ie when we plot a graph that should be separated correctly there is no misclassified points. then only we use hard margin s.v.m eqn

In Other case if we got some misclassified points in graph then we have a problem with that point we have to change or correctly classify the points

In Order to handle this problem. we have to change our eqn in to soft margin SVM eqn formulation where we treat this problem.



So if we observe the above graph point (1) is misclassified

point (2) - also misclassified

point (3) - Short distance from margin to misclassified point.

point (4) - Short " from " to correctly classified point

Point (5) - correctly classified point.

Here we solve the problem of Incorreally classified points (10)

for that we use signed distance in that we have actual value and predicted value's are presented. for point (1) we have actual is +ve predicted as (-ve) with some longer distance as shown in figure. shows the entire negative representation, ie misclassified point that can be represent with some ( $\xi_i$ ) Zeeta value. in point (1) we have value is Positive as Similarly if we do with all points.

P(2)  $\rightarrow$  Negative  $\rightarrow$  Misclassified  $\rightarrow$  ( $\xi_i$ ) positive

P(3)  $\rightarrow$  Negative  $\rightarrow$  Shorter distance  $\rightarrow$  misclassified  $\rightarrow$  ( $\xi_i$ ) positive

P(4)  $\rightarrow$  Positive  $\rightarrow$  Shorter distance  $\rightarrow$  Correctly classified  $\rightarrow$  ( $\xi_i$ ) Negative

P(5)  $\rightarrow$  Positive  $\rightarrow$  Correctly classified  $\rightarrow$  ( $\xi_i$ ) Negative.

from the above Summary we can say misclassified points having Positive ( $\xi_i$ ) value, for correctly classified points we have ( $\xi_i$ ) Negative value. that can be simply represent as -

when -  $\xi_i < 0 \Rightarrow$  (negative) Correctly classified. (lower than zero)

$\xi_i \geq 0 \Rightarrow$  (positive) Misclassified according to  $\pi_+$  &  $\pi_-$ .

$\xi_i \rightarrow$  called as Slack Variable.

Note :- if we increase zeeta value ( $\xi_i$ )  $\uparrow \rightarrow$  point is far away

in the incorrect direction ( $\pi_+$ ,  $\pi_-$ )

$y_i (w^T x_i + w_0) > 1 \Rightarrow \xi_i < 0$  (Correctly classified)

$y_i (w^T x_i + w_0) \leq 1 \Rightarrow \xi_i \geq 0$  (Misclassified)



So here Our aim is to maximize the margin. (11)

In Order to maximize margin we have to minimize function from that

$$\text{we can say } \text{Max}(\text{Margin}) = \text{Min} \frac{1}{\text{Margin}}.$$

from this eqn of Softmargin SVM will be -

$$\begin{aligned} \text{Min } & \left\{ \frac{1}{2} \|w\|^2 + c \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \right\} \\ \text{s.t. } & y_i (w^T x_i + w_0) \geq 1 - \xi_i \quad \forall_i \\ & \xi_i \geq 0. \end{aligned} \Rightarrow \text{Optimized Soft Margin SVM formulation}$$

Where  $c$  is hyperparameter of SVM.

from this eqn we can write formulation of Dual-form of S.M. SVM.

$$\begin{aligned} \text{Max } & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \underbrace{k(x_i, x_j)} \\ \text{s.t. } & 0 \leq \alpha_i \leq c \quad \& \quad \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

This term represent kernel trick

This entire thing is "SVM" with different formulation, In real time  
SVM is more cost ie most time Consuming algorithm. after K.N.N.  
The time Complexity is more for this algorithm, as Compare to K.N.N.  
SVM is less time Complexity.

K.N.N. > SVM > Logistic Regression

\* K.N.N. is also known as Lazy learner Algorithm.