# Support Vector Machine (SVM)
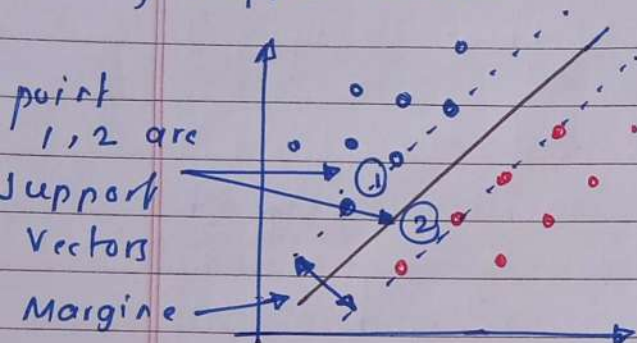
Basic about SVM.
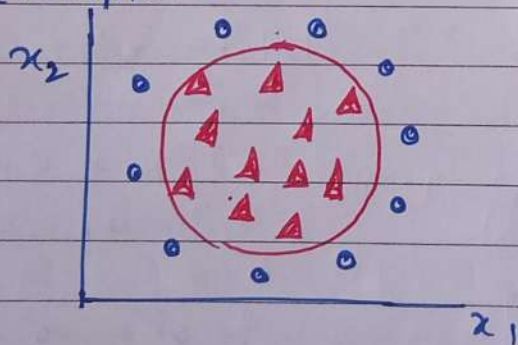
1) it is supervised ML model.
2) it can be used for both classification as well as regression but it is predominantly used for binary classification.
3) Hyperplane.
4) Support vectors

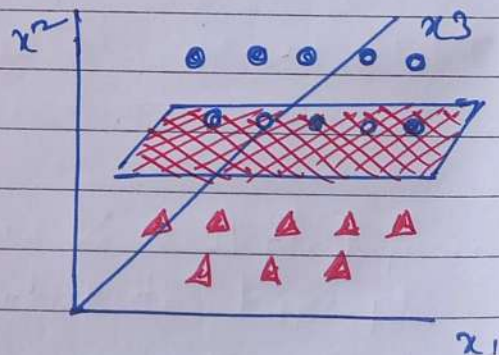point 1,2 are Support Vectors

Margine

position of hyperplane depends upon support Vectors.

- for 2D data it is easy to draw hyperplane but where data point can't seperated by line need to convert into 3D where we can seperate the datapoint by hyperplane Example :-

$x_2$

{ Here Not easy to seperate } by line

$x_1$

$x^2$

$x^3$

$2D \longrightarrow 3D$
and seperate by hyperplane

$x_1$

**Hyperplane :—** Hyperplane is line (in 2D) or plane that seperate the data point into two classes

<u>Support vectors</u> :- these are the datapoint which are nearest to hyperplane if these datapoints changes position of hyperplane changes.
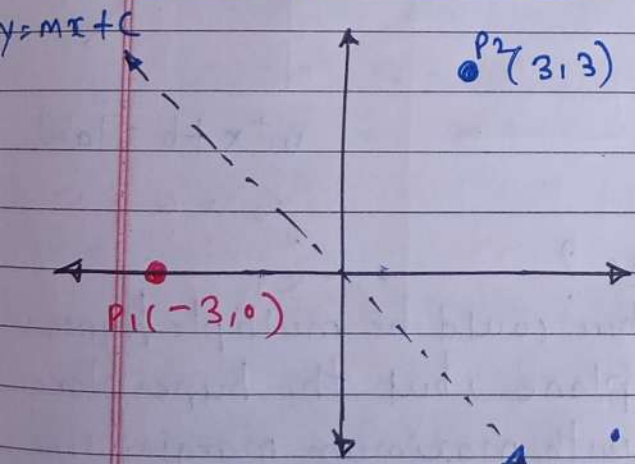
## <u>Advantages of SVM.</u>
1) works fine with smaller·dataset
2) works fine or efficiently where there is clear margine of seperation
3) works well with high dimensional data

## <u>Disadvantages</u>
1) Not suitable for large dataset as training time would take very large.
2) Not suitable for noiser (outlier) dataset with. overlapping classes.

## <u>Math Behind SVM.</u>

$y = mx + c$



$P_2(3,3)$

$P_1(-3,0)$

let slope and intercept of hyperplane is,

$m = -1$

$c = 0$ .. {since passing through origine}

• let parameters of hyperplane sare in W which is nothing but weight

$W \rightarrow (m,c) = (-1,0)$

• let multiply K or $P_1$ by transpose of w

$W^T x = \begin{bmatrix} -1 \\ 0 \end{bmatrix} \begin{bmatrix} -3 & 0 \end{bmatrix} = 3$  (positive)

[Note : why transpose ? → for matrix multiplication no of column of 1st Matrix must be equal ~~four~~ no. of rows of second matrix]

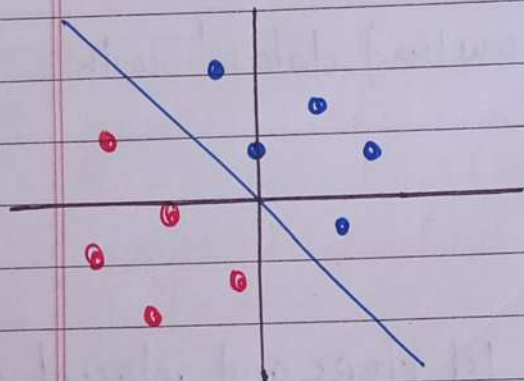- positive value indicates all the point of hyperplane will be positive class.

for $P_2(3,3)$

$$W^T x = \begin{bmatrix} -1 \\ 0 \end{bmatrix} [3,3]$$

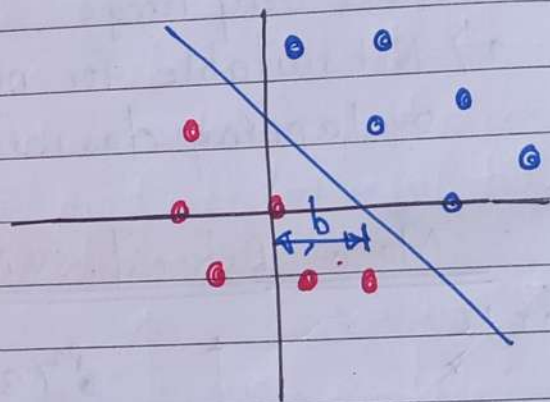$$= -3 \quad (\text{Negative})$$

Here for all the points which lie on the right side of hyperplane will be belong to negative class.

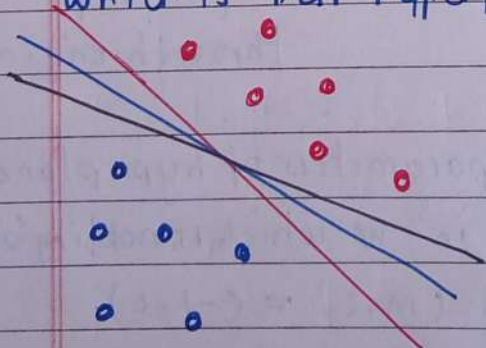But Not all the time hyperplane will pass through Origine.



$W^T x = label$

$W^T x + b = label.$
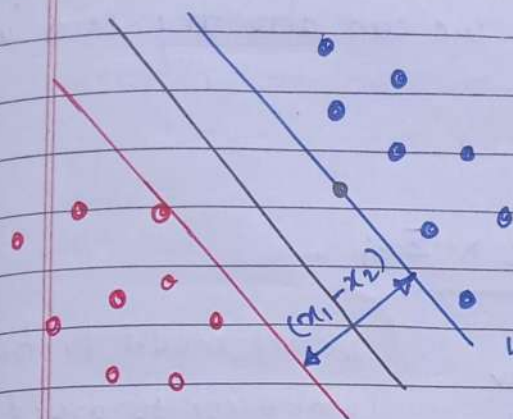
which is best hyperplane?



there could be multiple hyperplane, but the hyperplane with maximum margine size will be the best hyperplane.

$\longrightarrow$ optimization for maximum margine

$w^T x + b = label$

Equation of point or blue support vector. & its output value · any negative value.

$w^T x + b = -1$

$w^T x + b = 1 \Rightarrow$ this is equation of point or red support vector and its output value could be any positive value

to get margine let substract one from another.

$$w^T x_1 + b = 1$$
$$(-)\ w^T x_2 + b = -1$$
$$w^T (x_1 - x_2) = 2$$

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

divide both side by $\|w\|$

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

$$(x_1 - x_2) = \frac{2}{\|w\|} \longleftarrow \text{this is nothing but Magnitude of vector.}$$

and

$$Y_i = \begin{cases} -1 & w^T x_1 + b \leq -1 \\ 1 & w^T x_1 + b \geq 1 \end{cases} \quad (label)$$

So max $\left( \dfrac{2}{\|w\|} \right)$ such that.

$$Y_i = \begin{cases} -1 & w^T x_1 + b \leq -1 \\ 1 & w^T x_1 + b \geq 1 \end{cases}$$

instead of using. $\max_{\wedge}\left(\dfrac{2}{||w||}\right)$ we can also try Min which
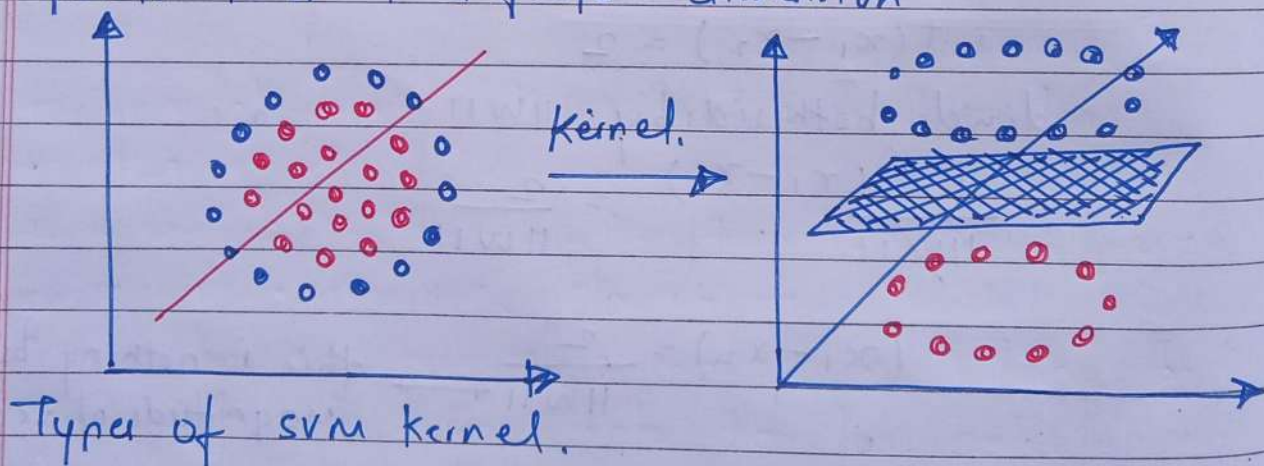
make better sense

$$\min\left(\dfrac{||w||}{2}\right) + \underbrace{c \times \Sigma \varepsilon_i}$$

c = Number of error

$\varepsilon_i$ = error magnitude

(we all model to train with some error to avoid overfitting (ie it will be good and train and will bad for test data)

<u>kernels in SVM</u> : Generally function of the kernel is to transform the training set of data so that non-linear decision surface can be transformed to a linear equation in higher number of dimension space it return the inner product between two points in standard feature dimension
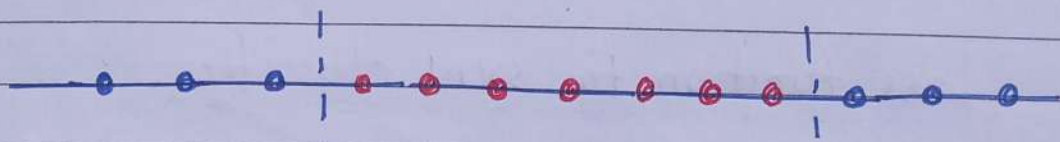


Kernel.

Types of SVM kernel.

1) Linear
2) polynomial
3) Radial Basis function. (rbf)
4) sigmoid.

Suppose feature $(x)$

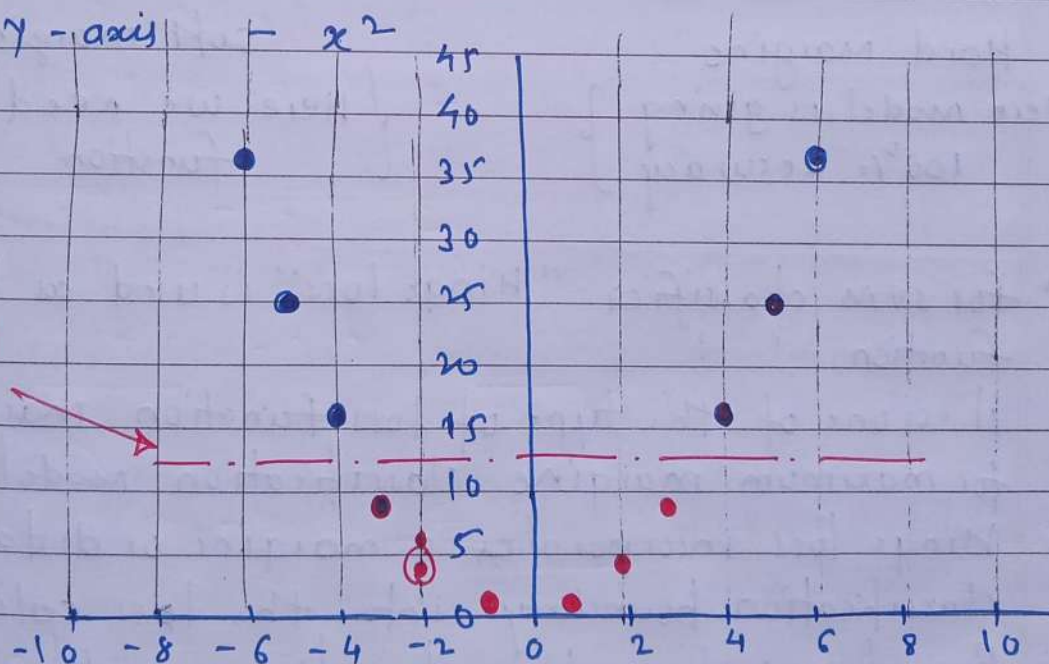| $x$ | $-6$ | $-5$ | $-4$ | $-3$ | $-2$ | $-1$ | $0$ | $1$ | $2$ | $3$ | $4$ | $5$ | $6$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $x^2$ | 36 | 25 | 16 | 9 | 4 | 1 | 0 | 1 | 4 | 9 | 16 | 25 | 36 |

if you try to plot $x$ on thin 1D line.

- we can see None of the line could seperate the two class pefectly.
- thal'r why we add another feature which is function of $x$    i.e $x^2$

    $x$ - axis — $x$
    $y$ - axis — $x^2$

Now this is sepuable data.

1) Linear kernel :- $k(x_1, x_2) = x_1^T \cdot x_2$
   { best suitable for having too many features }

2) polynomial kernel.
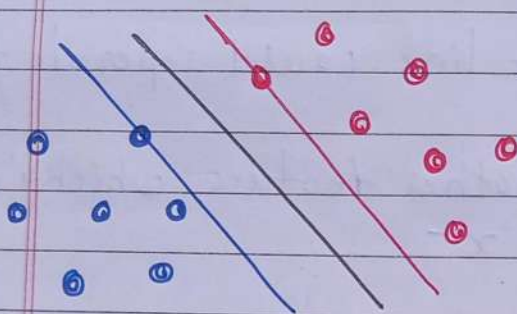   $k(x_1, x_2) = (x_1^T x_2 + r)^d$
                 degree

3.> radial basis function. (rbf kernel)

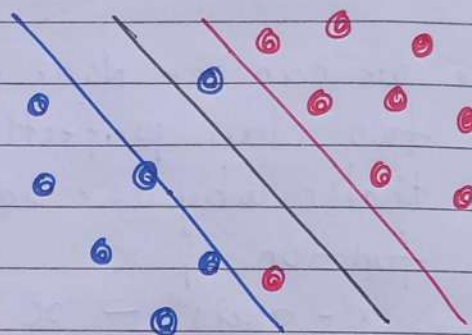$$k(x_1, x_2) = \exp(-\gamma \cdot ||x_1 - x_2||^2)$$

4'. Sigmoid function

$$k(x_1, x_2) = \tanh(\gamma \cdot x_1 + x_2 + r)$$

## Loss function for SVM classifier.



Hard Margine
{ Here model is giving
100% accuracy }
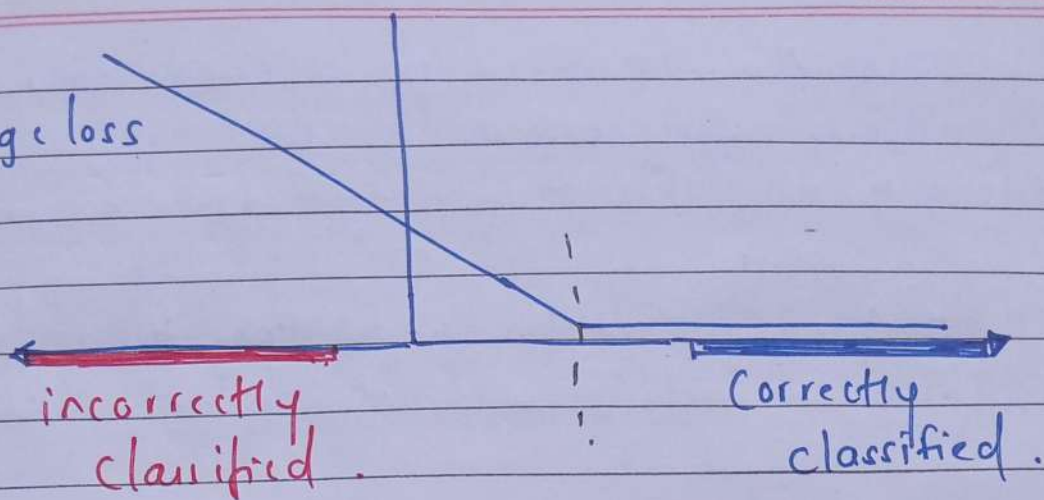
soft margine
{ Here we need loss
function }

- for SVM classifier "Hinge loss" is used as loss function.
- it is one of the type of loss function mainly used for maximum margine classification model.
- Hinge loss incorporates a margine or distance from classification boundary into the loss calculation Even if new observation classified correctly they can incure penalty if the margine from decision Boundary is not large enough.

$$L = \max(0, 1 - Y_i(w^Tx + b))$$

0 — for correct classification
1 — for wrong classification

hinge loss



incorrectly
classified.

Correctly
classified.

wrong

let tul for misclassification.

$Y_i = 1$, $\hat{Y_i} = -1$          $Y_i = -1$          $\hat{Y} = 1$

$L = (1 - 1(-1))$                     $L = (1 - (-1)(1))$

$= 1 + 1$                             $= 1 + 1$

$= 2$                                 $= 2$

{ both are high loss value }

Now tul for correct classification.

$Y_i = 1$    $\hat{Y_i} = 1$                     $Y_1 = -1$    $\hat{Y_1} = -1$

$(0 - (1)(1))$                                   $(0 - (-1)(-1))$

                                                 $0 - 1$

$0 - 1$                                          $-1$

$-1$

{ both are low loss value }