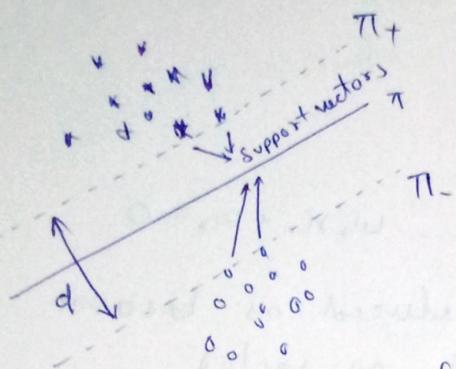


24-08-2023

# Support Vector Machine



find ~~best~~ hyperplane w/ margin ~~d~~  
maximize ~~w^T d~~  $\|w\|$   
margin =  $d$  = shortest distance b/w  
 $\pi_+$  and  $\pi_-$

margin-maximizing hyperplane

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

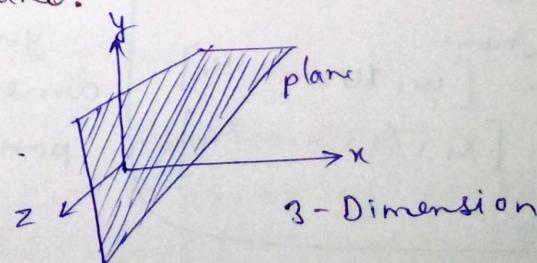
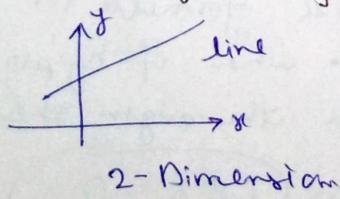
find the value of  $w$  and  $b$  where  $d$  is maximum.

$\pi$  = hyperplane

Advantages

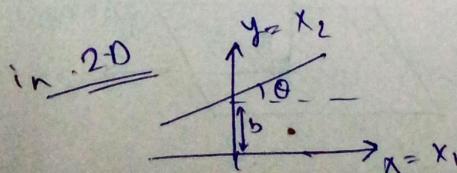
- SVM is very Robust to outliers (handle outliers)
- also work on non-linear data (using kernel)
- implied in both classification and regression both

equation of a hyperplane:



in 4D, 5D, 6D it's hyperplane.

- most of the dataset in ML are more than the 3 features



$$y = mx + b$$

general form

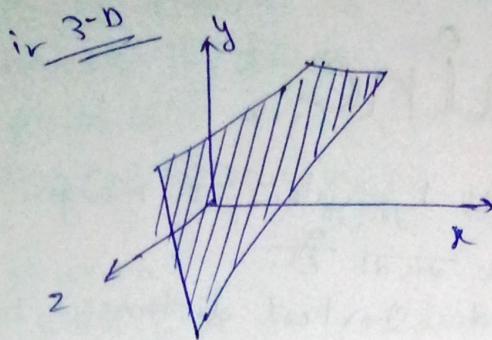
$$ax + by + c = 0$$

$$\text{where } y = -\frac{ax}{b} - \frac{c}{b} \text{ and } m = -\frac{a}{b}, b = -\frac{c}{b}$$

$$ax_1 + bx_2 + c = 0$$

$$\text{or } ax_1 + bx_2 + c = 0$$

$$\text{or } w_1x_1 + w_2x_2 + w_3 = 0$$



$$w_1x_1 + w_2x_2 + w_3x_3 + w_0 = 0$$

$$\text{in } n\text{-D} \quad w_1x_1 + w_2x_2 + w_3x_3 + \dots + w_nx_n + w_0 = 0$$

In the above equations can be reduced as because the given points can be used as vectors.

$$\underbrace{w \cdot x + w_0 = 0}$$

$$[w_1, w_2, \dots, w_n] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} + w_0 = 0$$

$$w^T x + w_0 = 0$$

Generalize formula for any dimension when if hyperplane pass from the origin then  $w_0 = 0$

$$\text{then } w^T x = 0$$

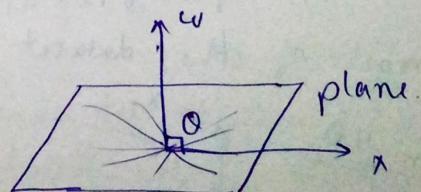
$$w^T x = w \cdot x = \|w\| \|x\| \cos \theta = 0$$

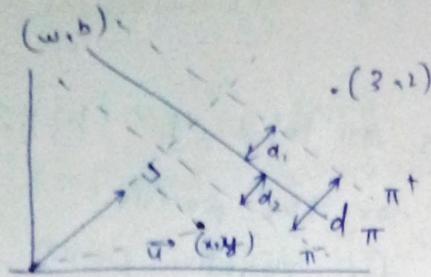
$$\text{if } \theta = 90^\circ$$

it means by diagram

$w$  is  $\perp$  to the plane

E&QII





$$w \cdot u \geq c$$

$$w \cdot u - c \geq 0$$

$$w \cdot u + b \geq 0$$

decision rule.

example:  $2x_1 + 3x_2 + 3 = 0$

$$w_1 x_1 + w_2 x_2 + 3 = 0$$

$$\vec{w} = (2, 3) \quad \vec{u} = (3, 2)$$

$$2(3) + 3(2) + 3 > 0$$

then  $(3, 2)$  is lies in +ve side

similarly for  $\vec{u} = (-2, 0)$  point is lies in -ve side

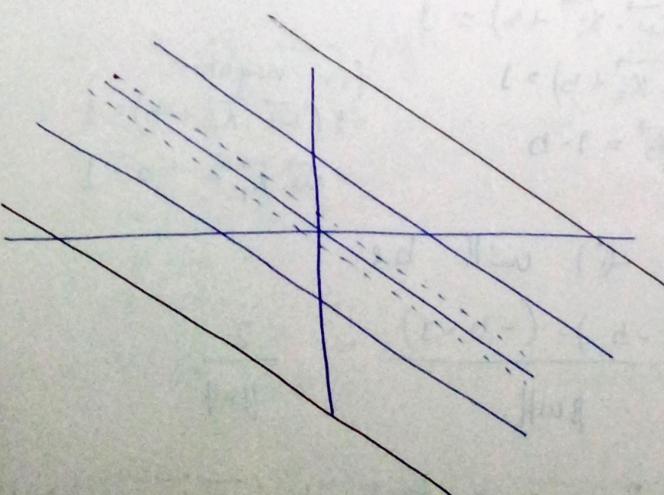
for any  $x_i$

$$\hat{y} = \begin{cases} +1 & \text{if } \vec{w} \cdot \vec{x}_i + b \geq 0 \\ -1 & \text{if } \vec{w} \cdot \vec{x}_i + b < 0 \end{cases}$$

assume that

$$\text{for } \pi_+: \vec{w}^T \vec{x} + b = 1$$

$$\text{for } \pi_-: \vec{w}^T \vec{x} + b = -1$$



let us example in 2-D

$$\begin{cases} 2x_1 + 3x_2 + 3 = 1 \\ 2x_1 + 3x_2 + 3 = 0 \\ 2x_1 + 3x_2 + 3 = -1 \end{cases}$$

$$\begin{cases} 2x_1 + 3x_2 + 3 = 1 \\ 2x_1 + 3x_2 + 3 = 0 \\ 2x_1 + 3x_2 + 3 = -1 \end{cases}$$

$$\begin{cases} 0.2x_1 + 0.3x_2 + 0.3 = 1 \\ 0.2x_1 + 0.3x_2 + 0.3 = 0 \\ 0.2x_1 + 0.3x_2 + 0.3 = -1 \end{cases}$$

Defining constraints of SVM

$$\begin{cases} \vec{w} \cdot \vec{x}_i + b \geq 1 \\ \vec{w} \cdot \vec{x}_i + b \leq -1 \end{cases}$$

$\vec{x}_i$  are and are points and label  $y_i = +1, y_i = -1$

$y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1 \rightarrow$  for +ve points

$y_i(\vec{w} \cdot \vec{x}_i + b) \geq -1 \rightarrow$  for negative points

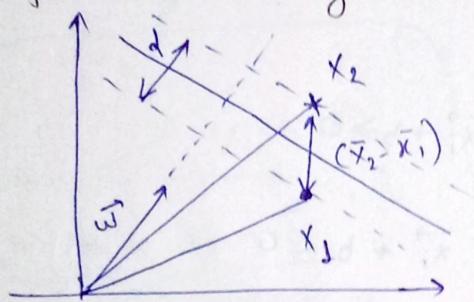
our final constraint is

$$y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1$$

for support vectors

$$y_i(\vec{w} \cdot \vec{x}_i + b) = 1$$

find max margin in SVM:



Now unit vector of  $w = \frac{\vec{w}}{\|w\|}$

distance vector b/w  $x_1, x_2 = (\vec{x}_2 - \vec{x}_1)$

$$d_2 = (\vec{x}_2 - \vec{x}_1) \cdot \frac{\vec{w}}{\|w\|}$$

$$d_1 = \frac{\vec{x}_2 \cdot w - \vec{w} \cdot \vec{x}_1}{\|w\|} \quad (1)$$

$$\therefore y_i(\vec{w} \cdot \vec{x}_i + b) = 1$$

for +ve  $\vec{w} \cdot \vec{x}_i + b \geq 1$

$$\vec{w} \cdot \vec{x}_i = 1 - b$$

for negative  $-\vec{w} \cdot \vec{x}_i + b = 1$

$$\vec{w} \cdot \vec{x}_i = -b - 1$$

hence equation (1) will be

$$d_2 = \frac{p(1-b) - (-b-1)}{\|w\|} = \frac{2}{\|w\|}$$

Hence  $\text{argmax } (\vec{w}^*, b^*) \frac{2}{\|w\|}$  such that  $y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1$

always true for every point given.

all the above calculation and finding the maximum margin we called it Hard margin SVM

Soft Margin SVM:  
constraint optimization problem

$$\max f(x) \Leftrightarrow \min \frac{1}{f(x)}$$

Hence.

$$C \propto \frac{1}{\lambda}$$

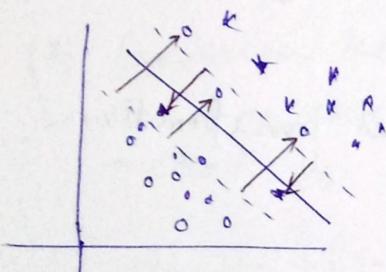
$$d = \underset{(w^*, b^*)}{\operatorname{arg\,min}} \left[ \frac{\|w\|}{2} + C \sum_{i=1}^n \xi_i \right] \rightarrow \text{Hinge loss}$$

$\xrightarrow{\text{regularization}}$  Here, the above formula is regularization

$C$  = hyper parameter

$\xi$  = zeta =

if all points are correctly classified. then  $\sum_{i=1}^n \xi_i = 0$

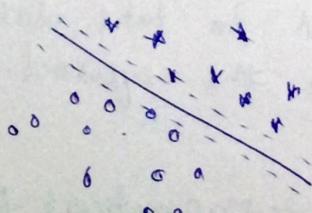


$\sum_{i=1}^n \xi_i$  is the sum of all the distance of incorrectly classified point

Home in SVM

Error = Margin error + classification error

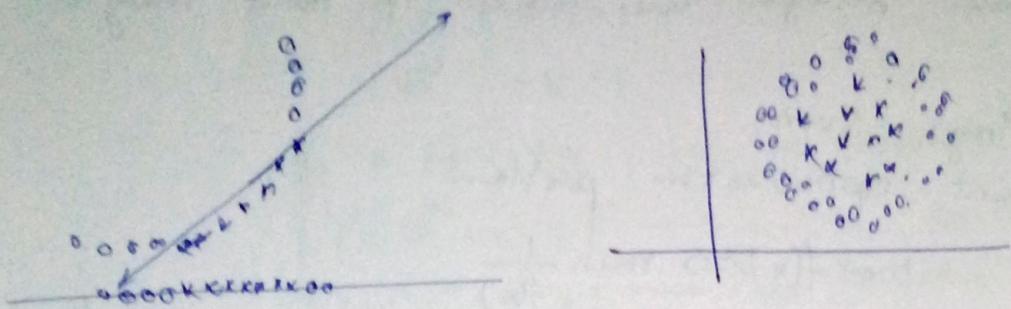
if  $C$  is very high.



if  $C$  is very low



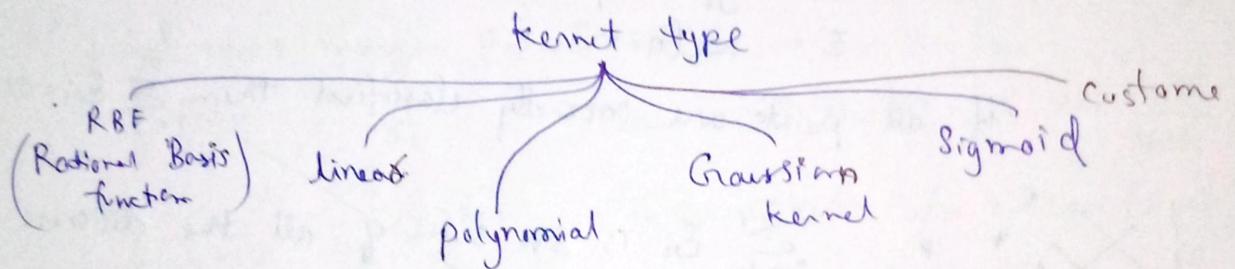
## Kernel trick in SVM



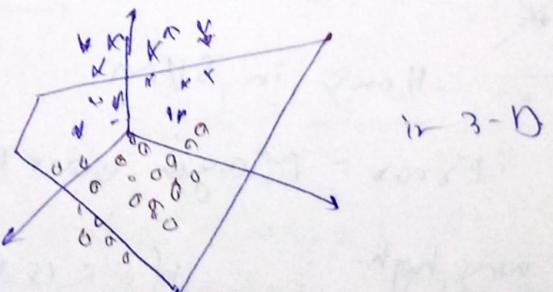
Non-linear

data

(low-dimension  $\rightarrow$   $f(x)$ )  $\rightarrow$  higher dimension.  
kernel



- The transformation is called kernel transformation.



- Kernel function is a method used to take data as input and transform it into the required form of processing data.

- Kernel is used due to a set of mathematical functions used in SVM providing the window to manipulate the data.

- So kernel function generally transforms the training set of data so that a non-linear

decision surface is able to transform to a linear equation in higher number of dimension spaces.

- Basically  $\phi H$  returns the inner product b/w two point in a standard feature dimension.

standard kernel function equation

$$k(\bar{x}) = \begin{cases} 1 & \text{if } \|\bar{x}\| = 1 \\ 0 & \text{otherwise} \end{cases}$$

(i) Gaussian kernel:

- it is used to perform transformation when there is no prior knowledge

$$k(x, y) = e^{-\frac{\|(x-y)\|^2}{2\sigma^2}}$$

(ii) Gaussian Kernel Radial Basis Function: (RBF):

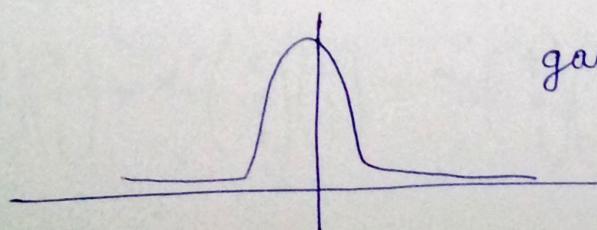
- Same as above kernel function, adding radial basis method to improve the transformation.

$$k(x, y) = e^{-\frac{\|(x-y)\|^2}{\gamma}}$$

$$K(x_1, x_2) + K(x_1, x_3) \quad (\text{simplified formula})$$

$$K(x_1, x_1) + K(x_1, x_2) > 0 \quad (\text{Green})$$

$$K(x_1, x_1) + K(x_1, x_2) = 0 \quad (\text{Red})$$



gaussian kernel graph.

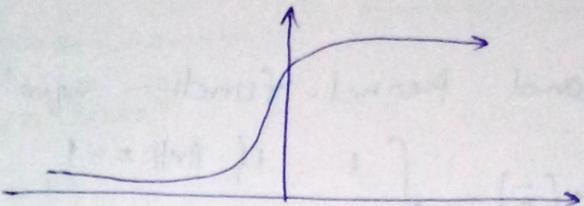
code:  
from sklearn.svm import SVC  
clf = SVC(kernel='rbf', random\_state=0)  
clf.fit(X\_train, y\_train)

\* if  $\gamma$  is very high then the model overfits the training data

### (iii) Sigmoid kernel.

- This function is equal to a two-layer perceptron model of the neural network, which is used as an activation function for artificial neurons.

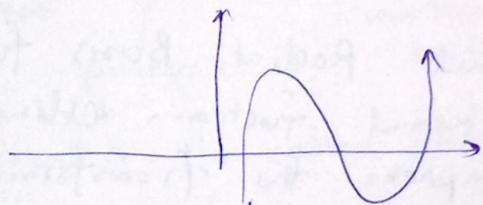
$$K(x, y) = \tanh(x \cdot x^T y + \gamma)$$



### (iv) polynomial kernel:

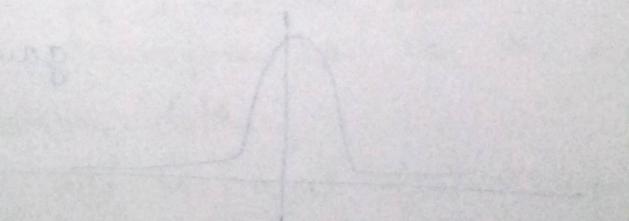
- It represents the similarity of vectors in the training set of data in a feature space over polynomials of the original variables used in the kernel.

$$K(x, y) = \tanh(x \cdot x^T y + \gamma)^d, \gamma > 0$$

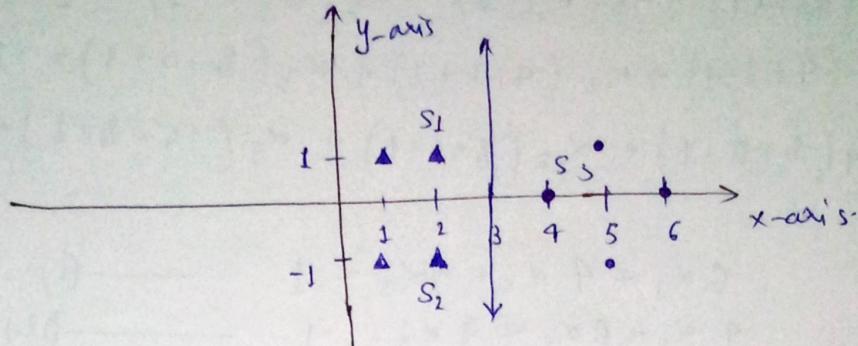


### (v) linear kernel: used when data is linearly separable.

$$\begin{aligned} \text{(new)} &= c(\text{old}) + \alpha(x_{\text{new}}) \\ \text{(last)} &= c - (\alpha(x_{\text{old}})) + \alpha(x_{\text{new}}) \end{aligned}$$



Question Find the ~~optimal~~ hyper-plane for following data points: negatively labelled data points.  $(1, 2), (2, 1), (2, -1) \text{ and } (2, -1)$  and positively labelled data points.  $(4, 0), (5, 1), (5, -1) \text{ and } (6, 0)$



Step-1 find the support vectors

$$S_1 = \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} \quad S_2 = \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} \quad S_3 = \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix}$$

Step-2. finding augmented vector

$$S_1 = \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} \quad S_2 = \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} \quad S_3 = \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix}$$

Step 3 finding expressions.

$$\alpha_1 \bar{S}_1 \cdot \bar{S}_1 + \alpha_2 \bar{S}_2 \cdot \bar{S}_1 + \alpha_3 \bar{S}_3 \cdot \bar{S}_1 = -1$$

$$\alpha_1 \bar{S}_1 \cdot \bar{S}_2 + \alpha_1 \bar{S}_2 \cdot \bar{S}_2 + \alpha_3 \bar{S}_3 \cdot \bar{S}_2 = -1$$

$$\alpha_1 \bar{S}_1 \cdot \bar{S}_3 + \alpha_2 \bar{S}_2 \cdot \bar{S}_3 + \alpha_3 \bar{S}_3 \cdot \bar{S}_3 = 1$$

put the value of  $\bar{S}_i$  in above equation

$$\alpha_1 \begin{pmatrix} ? \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} = -1$$

$$\alpha_1 \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} = -1$$

$$\alpha_1 \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 9 \\ 0 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} = 1$$

Now,

$$\alpha_1(9+1+1) + \alpha_2(-1+1) + \alpha_3(8+0+1) = -1$$

$$\alpha_1(9-1+1) + \alpha_2(4+1+1) + \alpha_3(8+0+1) = -1$$

$$\alpha_1(8+0+1) + \alpha_2(8+0+1) + \alpha_3(16+8+1) = 1$$

$$6\alpha_1 + 4\alpha_2 + 9\alpha_3 = -1 \quad \text{--- (i)}$$

$$4\alpha_1 + 6\alpha_2 + 9\alpha_3 = -1 \quad \text{--- (ii)}$$

$$9\alpha_1 + 9\alpha_2 + 17\alpha_3 = 1 \quad \text{--- (iii)}$$

By solving the above equation we get.

$$\alpha_1 = \frac{-13}{4} = -3.25; \alpha_2 = \frac{-13}{4} = -3.25; \alpha_3 = \frac{7}{2} = 3.5$$

Step-4 Calculate the weight vector:

$$w_2 \sum \vec{\alpha_i s_i}$$

$$w = \alpha_1 \vec{s}_1 + \alpha_2 \vec{s}_2 + \alpha_3 \vec{s}_3$$

$$(w, b) = -3.25 \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + -3.25 \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} + 3.5 \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix}$$

$$(w, b) = \begin{pmatrix} -6.5 + 6.5 + 14 \\ -3.25 + 3.25 + 0 \\ -3.25 + -3.25 + 3.5 \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \\ -3 \end{pmatrix}$$

$$(w, b) = \begin{pmatrix} 1 \\ 0 \\ -3 \end{pmatrix}$$

Step-5 Hyperplane equation

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

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

$$\cancel{w = } \quad w = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad b = -3$$

Hence hyper plane line will pass through at  
 $x=3$  which will be parallel to y-axis.

1. What are kernels in SVM? Can you list some popular SVM kernels?

Description:

Kernels in SVM are functions that transform data into a higher-dimensional space, making it easier to find a hyperplane (decision boundary) that separates data points. These functions enable SVMs to handle non-linear data by projecting it into higher dimensions, essential for capturing complex patterns.

Advantages:

Kernels are versatile and effective in capturing complex patterns, allowing SVMs to handle non-linear data. They transform data to enable the separation of complex patterns.

Disadvantage:

Choosing the right kernel and its parameters can be challenging, and training with complex kernel can be computationally expensive. Proper kernel selection is crucial.

Numerical example:

In a binary classification problem, an RBF kernel projects data into a higher-dimensional space, making it separable by a hyperplane and allowing the SVM to find a non-linear decision boundary.

Popular SVM kernels:

- (i) Linear kernel: computes the dot product in the original space, suitable for linearly separable data.

- (ii) Polynomial kernel: Raises the dot product to a power, introducing non-linearity and is suitable for data with polynomial boundaries.
- (iii) Radial Basis Function (RBF) kernel:  
uses a Gaussian-like function to capture complex non-linear relationships, and is widely used and versatile.
- (iv) Sigmoid kernel:  
Applies the hyperbolic tangent function to the dot product, useful for data with sigmoid-shaped boundaries.
- (v) Custom kernels:  
Custom kernels can be defined based on domain knowledge or problem-specific features.