# Support Vector Machine
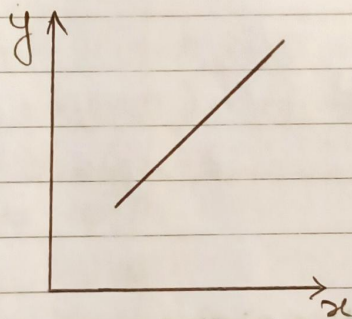
It can solve both classification and Regression problem.

1. classification → SVC (support Vector classifier)

2. Regression → SVR (support Vector Regressor)

## some basics:



Equation of line:

$$y = mx + c \quad OR$$
$$y = \beta_0 + \beta_1 x \quad OR$$

$$ax + by + c = 0$$

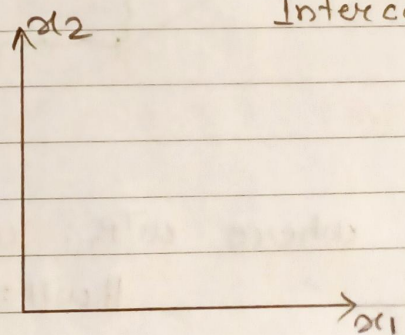$$\therefore \quad y = \underbrace{\frac{-a}{b}}_{\text{coefficient}} x - \underbrace{\frac{c}{b}}_{\text{Intercept}}$$

$$ax_1 + bx_2 + c = 0$$
$$w_1 x_1 + w_2 x_2 + b = 0$$
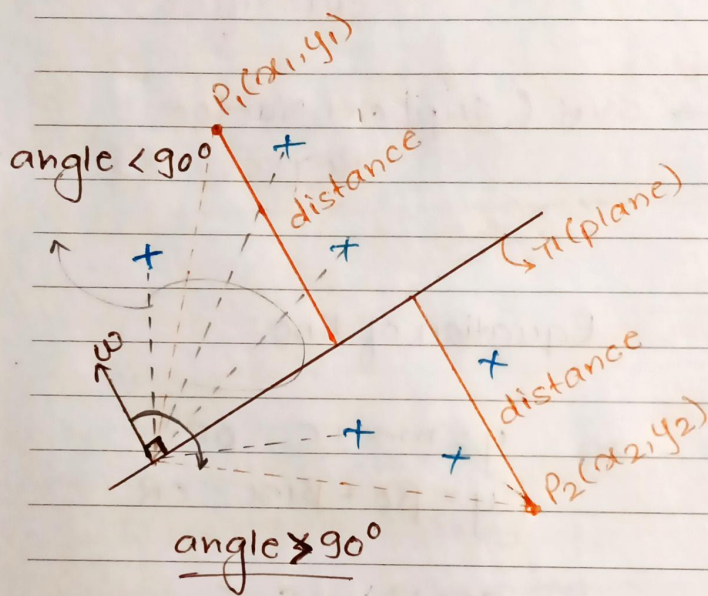
$$\therefore \quad w^T x + b = 0$$

If line passes through origin $\therefore$ $w^T x = 0$

matrix multiplication: $\begin{bmatrix} \omega_1 \\ \omega_2 \end{bmatrix} \begin{bmatrix} x_1 & x_2 \end{bmatrix}$

$$\omega^T x + b = 0 \qquad (\omega^T x : \omega \text{ Transpose } x)$$

Equation of line passing through origin is:

$$\boxed{\omega^T x = 0}$$

angle < 90°

$P_1(x_1, y_1)$

distance

$\pi$ (plane)

$\omega$

angle > 90°

distance

$P_2(x_2, y_2)$

we have to find the distance of point from the plane.

$(\pi = \text{line in 2D}$
$\text{and plane in 3D})$

Distance of a point to the plane,

$$\boxed{\text{distance } (d) = \dfrac{\omega^T P_1}{\|\omega\|}}$$

$$= \|\omega\| \cdot \|P\| \cdot \cos\theta$$

where $\omega^T P_1$ : $\omega$ Transpose $P_1$,  $\omega$ : vector,
$\|\omega\|$ : magnitude of $\omega$

**unit vector:** A vector which has a magnitude of 1 is basically called unit vector.

Eg.



Now,

$$d = \sqrt{3^2 + 4^2}$$
$$= \sqrt{25}$$
$$\therefore \quad d = 5$$

where vector, $\hat{d} = \dfrac{d}{\boxed{\lVert d \rVert}} \rightarrow$ magnitude

$$\left(3/5, 4/5\right) = d = \sqrt{\left(3/5\right)^2 + \left(4/5\right)^2} = \sqrt{25/25}$$
$$= 1.$$

$\therefore$ unit vector is a way to get focused on direction not on magnitude.

**upward vector** →

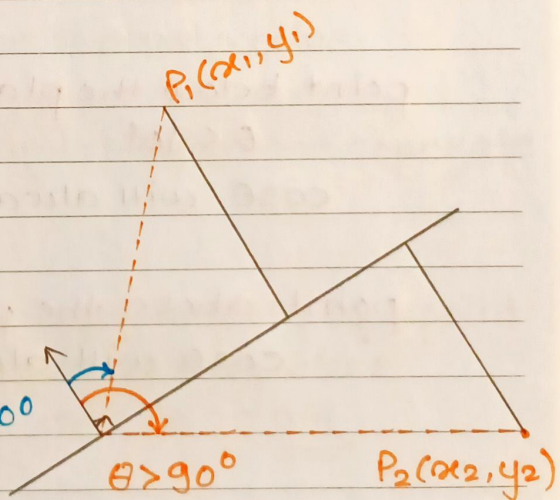$$d = \dfrac{w^T P_1}{\lVert w \rVert}$$
$$d = \lVert w \rVert \cdot \lVert P_1 \rVert \cdot \cos\theta$$



**point above the plane**
as $\theta < 90°$
$\theta < 90°$
$\cos\theta$ will always be +ve.
$\theta > 90°$

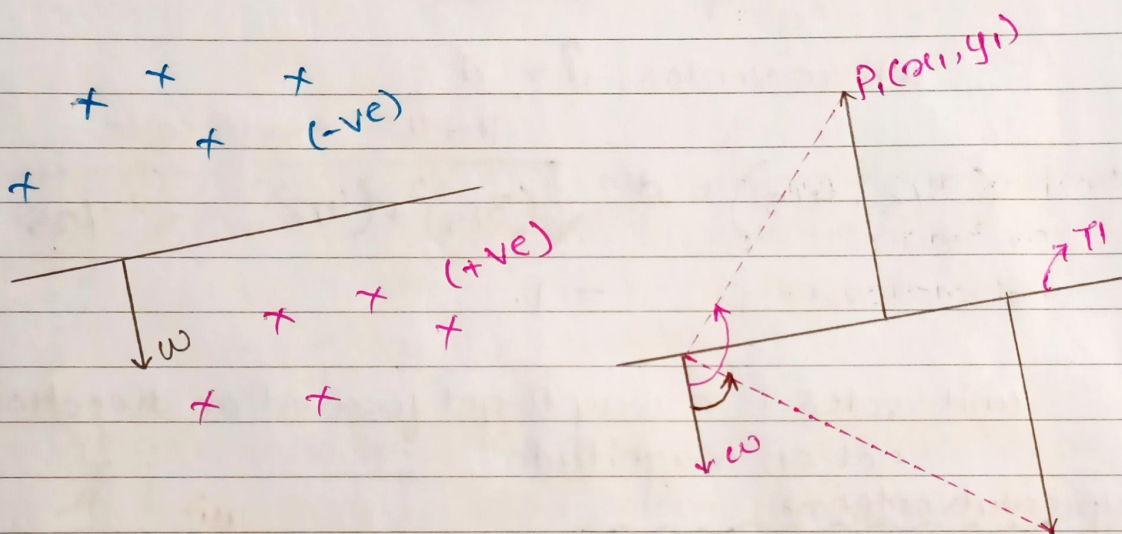$P_1(x_1, y_1)$
$P_2(x_2, y_2)$

**point below the plane,** as $\theta > 90°$
$\cos\theta$ will always be -ve.

- If any point falling above the plane, then θ must be less than 90° (+value)

- If any point falling below the plane, then θ must be greater than 90°. (-ve value)

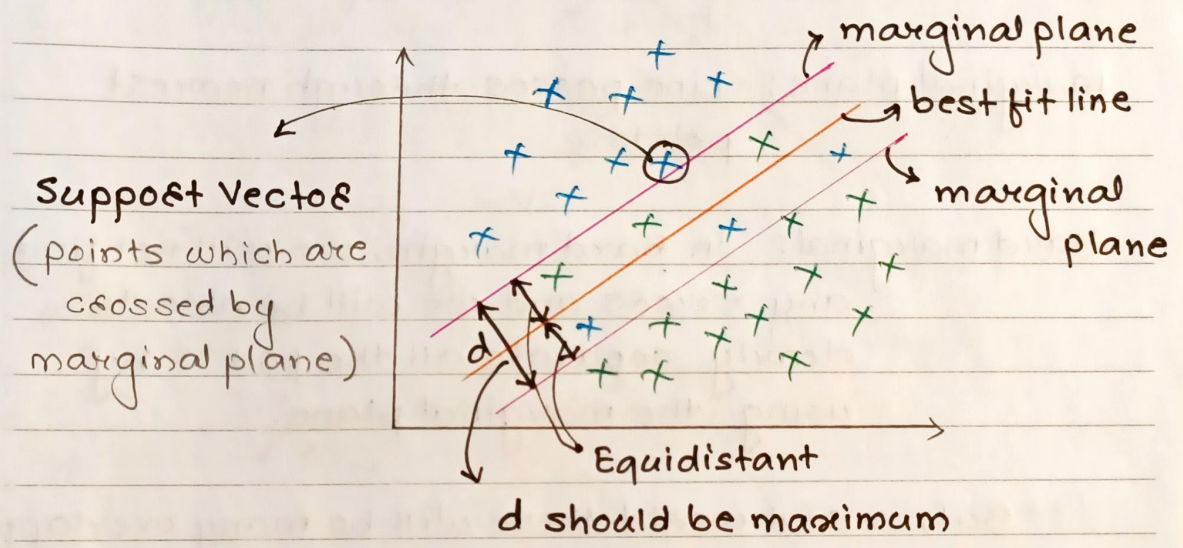## Downward vector:



point below the plane,
$\theta < 90°$
cosθ will always be '+ve'.

point above the plane, $\theta > 90°$
cosθ will always be '-ve'.

# Geometric Intuition Behind Support Vector Machine

## Support Vector classifier (SVC)



**Support Vector**
(points which are crossed by marginal plane)

marginal plane
best fit line
marginal plane

Equidistant

d should be maximum

* you can have more than one support vector.

There will be many possible hyperplanes that separate different classes.

we have learnt in LR that the probability of a point belonging to any class given at very close to the hyperplane will be close to 0.5.

So, we want a hyperplane that seperates (+ve) pts and (-ve) pts as far away as possible.

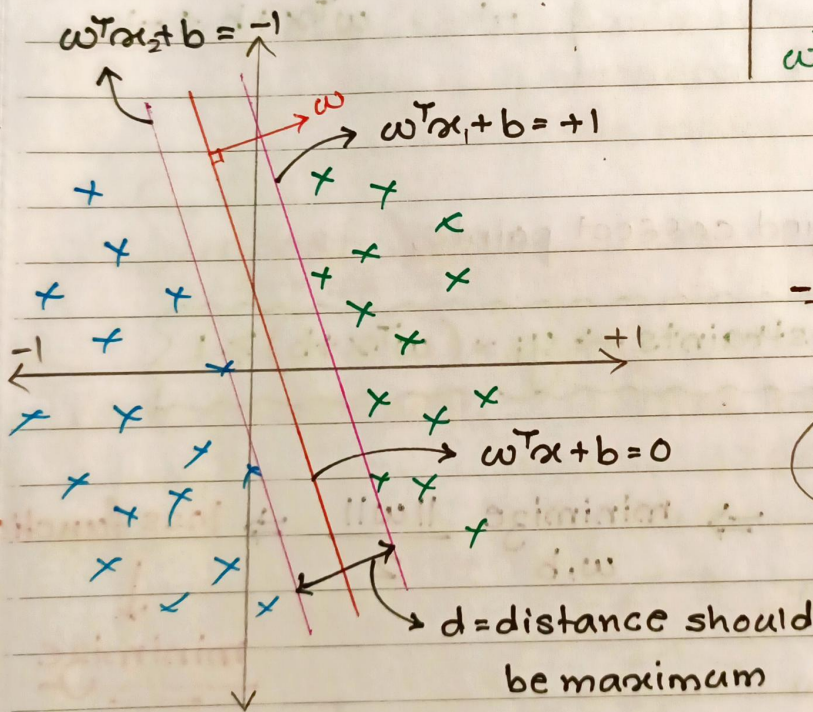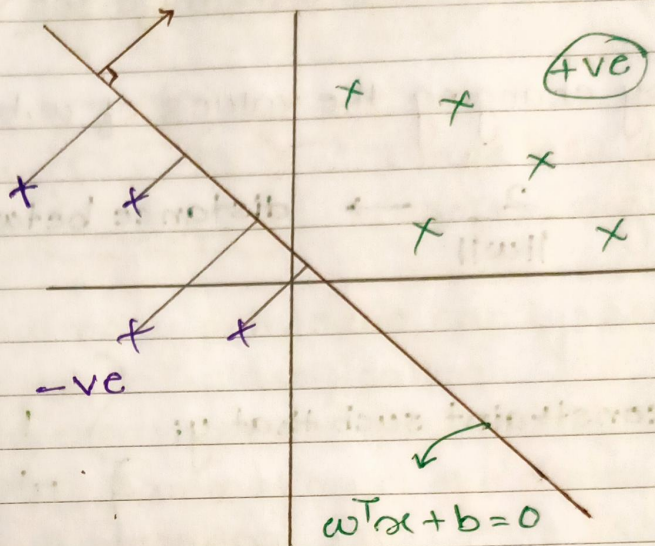↳ key idea of SVM. such hyperplane is called margin-maximizing plane.

Hard marginal                    Soft marginal.

**Marginal plane:** Line passes through nearest points.

**Hard marginal:** In hard margin, we will not find any errors and we will be able to clearly separate all the points by using the marginal plane.

But, in real world there will be many overlapping with many errors, so marginal plane lines will be called soft margin.

→ marginal plane should be equidistance from best fit line.

## SVM Mathematical Intuition:

(+ve)

+ve

$\omega^T x + b = 0$

$\omega^T x_2 + b = -1$

$\omega^T x_1 + b = +1$

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

$\omega^T x + b = 0$

$\therefore$ unit vector of $\omega$,

$\dfrac{\omega^T (x_1 - x_2)}{||\omega||}$

$= \dfrac{2}{||\omega||}$

d = distance should be maximum

$\therefore$ distance $(d) = \dfrac{\omega^T(x_1 - x_2)}{||\omega||} = \dfrac{2}{||\omega||}$

## Cost Function:

we have to maximize the value of $\frac{2}{\|w\|}$
by changing the values of $w, b$.

$\frac{2}{\|w\|}$ ⟶ distance between marginal plane.

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

~~conditions~~

For all classified correct points,

constraints ⟶ $y_i \times (w^T x + b) \geq 1$

$\underset{w,b}{\text{maximize}} \ \frac{2}{\|w\|}$ ⟹ $\underset{w,b}{\text{minimize}} \ \frac{\|w\|}{2}$ ⟹ loss function

↓

minimize

✳ loss function focus on minimization.

## cost function:

minimize $\dfrac{\|w\|}{2}$ by changing $w, b$.

$$\min \dfrac{\|w\|}{2} + C_i \sum_{i=1}^{3} \xi_i \longrightarrow \text{Hinge loss for soft margin.}$$

where, $C_i$ : How many points we can ignore for mis-classification.
$\hookrightarrow$ Hyperparameter

$\xi_i$ (Eta): Summation of the distance of incorrect data points from the marginal plane.

## Support Vector Regressor:

Problem Statement: Based on the size of the house, we have to predict price of the house.



price

$w^T x + b + \epsilon$

$w^T x + b$

$w^T x + b - \epsilon$

$\epsilon$ {
$\epsilon$ {

$\longrightarrow$ Size

$\{ \epsilon:$ epsilon $\rightarrow$ marginal error $\}$

## cost Function:

$$\text{minimize}_{w,b} \quad \frac{\|w\|}{2} + \left\{ C_i \sum_{i=1}^{w} \varepsilon_i \right\} \rightarrow \text{Hinge loss}$$

MAE

$$\text{constraint:} \quad |y_i - w_1 x_i| \leq \epsilon + \varepsilon_i \qquad \rightarrow \text{Eta qi}$$

Truth point     psedicted point     epcilon

$\epsilon$ : margin of erroε (to decide oεiginal plane)

$\varepsilon_i$ : erroε above the margin.



## Hyperparameter:

→ keep adjusting $\epsilon$ to get best margin

→ we can't say incoεεect point in εegεessoε.

Price (y-axis), Size (x-axis)

Lines labeled:
$w^T x + b + \epsilon$
$w^T x + b$
$w^T x + b - \epsilon$

predicted.

In Regressor, no complete incorrect value, because it will be continuous value.
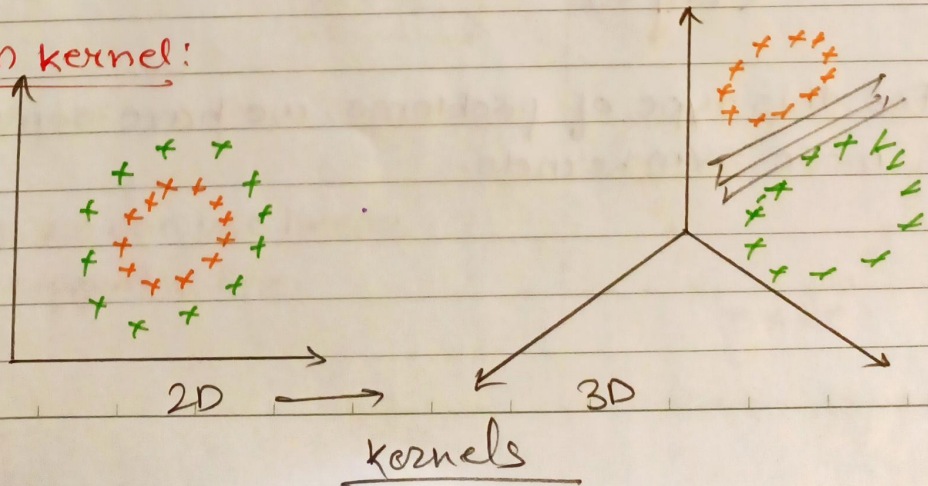
Q.

Is svm impacted by the outliers?
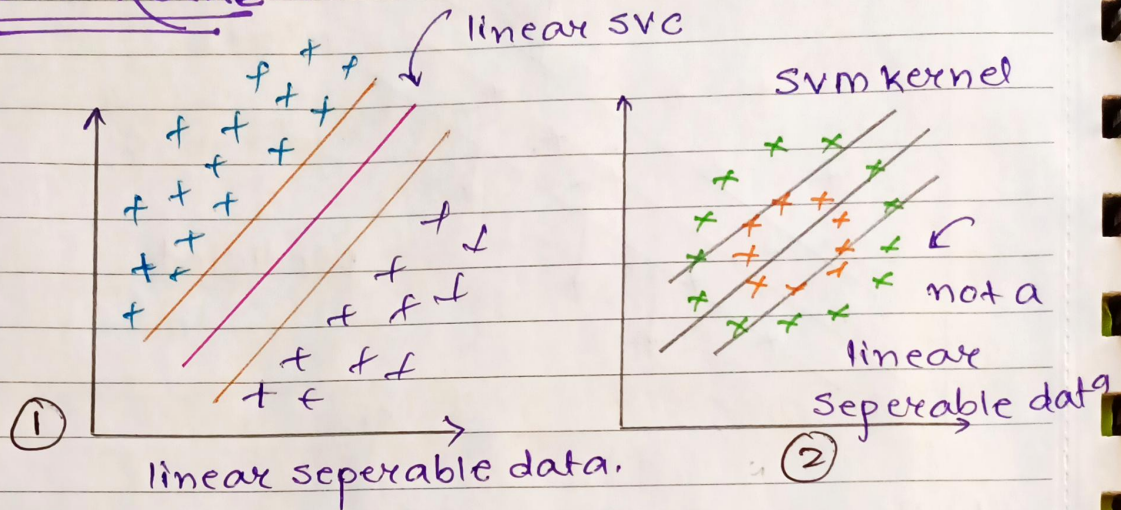Yes, svm is impacted by the outliers.

Does Standardization is need in svm?
Yes, we need to perform Normalization and standardization.

SVM kernel!



2D ⟶ 3D

Kernels

# SVM Kernels:



linear SVC

SVM kernel

① linear seperable data.

② not a linear seperable data

- when we create this (①) type of best fit line and marginal plane, we are actually solving the linear seperable data.

  - → called as Linear SVC. (Fig ①)

- If data is not a linear seperable data, you will not be able to create best fit line and not able to create a marginal plane even though we create it, the accuracy will be very low.
    (Fig 2)

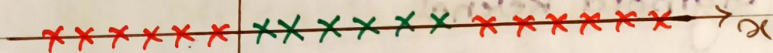- For this type of problems, we have some more SVM kernels.

**what does SVM kernels do?**

→ The main aim is to apply some transformation technique. (some mathematical formula) on the dataset.

This transformation increases the dimension of the data.

(mathematical formula)

SVM kernels → Transformation → Increasing the dimension of the data

linear seperable line
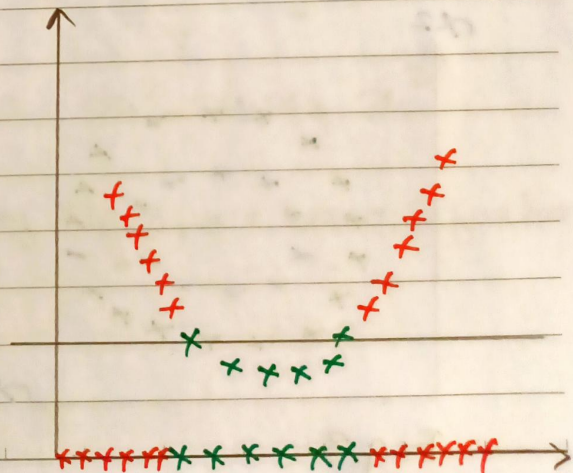→ will give error

×××××× ×××××× ×××××× → $x$

∴ we will transform the data from 1D → 2D.

$$y = x^2$$

After →

Now, we can use linear seperable line.

$$y = x^2$$

so if $x = -7$  $y = 49$

$x = -3$  $y = 9$  and so on.
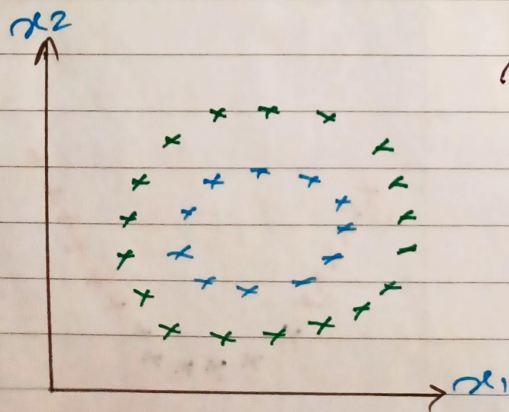
what is the advantage of doing this transformation?

→ after transformation, we can apply linear SVM or SVC.

  * when we ~~divide~~ convert 1D → 2D then we can divide all the points using single line which is called Linear SVC.

## Types of SVM kernel :

1. Polynomial kernel
2. RBF kernel
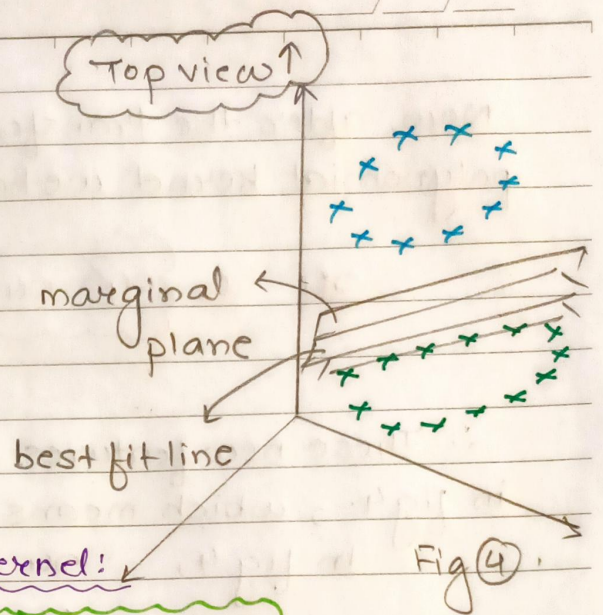3. Sigmoid kernel

## 1. Polynomial kernel



→ not seperable with best fit line.

so, we need to convert 2D → 3D.

Fig 3.

Our main aim was to
increase dimension,
$$2D \rightarrow 3D$$
so, hyperplane is
created.



Top view ↑

marginal
plane

best fit line

Fig ④

**Formula for Polynomial kernel:**

$$f(x_1, x_2) = \left( x_1^T x_2 + 1 \right)^d$$

d = dimension

If we are converting $2D \rightarrow 3D$, the value of d=3.

$$\therefore \quad x_1^T x_2 = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot \begin{bmatrix} x_1 & x_2 \end{bmatrix}$$

$$= \begin{bmatrix} x_1^2 & x_1 \cdot x_2 \\ x_1 \cdot x_2 & x_2^2 \end{bmatrix}$$

3 unique values: $x_1^2, x_1 \cdot x_2, x_2^2$

Now, initially at the time of Fig 3, we have 3 features
$$x_1, x_2, y$$

Now, after the transformation / formula of polynomial kernel we have 6 features

$$x_1 \quad x_2 \quad \overbrace{x_1^2 \quad x_1 . x_2 \quad x_2^2} \quad y$$

∴ These new features can be plotted as the 3D in fig 4. , which means that

In fig 4,  
$x_1$ will be $x_1^2$  
$x_2$ will be $x_2^2$  
$z$ will be $x_1 . x_2$

and once we have all these points, we will be able to clearly seperate the points.

→ use polynomial kernel, to get better accuracy.

② <u>Radial Basis Function Kernel (RBF kernel)</u>

$$K(\vec{x}, \vec{li}) = e^{- \dfrac{\| \vec{x} - \vec{li} \|}{2 6^2}}$$

↳ hyperparameter

③ <u>Sigmoid Kernel</u>

It can be used as the proxy for neural networks.

$$K(x, x_i) = \tanh\left(6 X^T X_i + \gamma\right)$$