

15/2/23

SVM (SUPPORT VECTOR MACHINE)

(most powerful & imp algo in ML)

Topics -

- Formulating Objective
- Objective is a Constraint Convex Function
- Handling Outliers
- Weight & Bias Update Rule
- Hinge loss function
- Handling non-linearly separable data
- Kernel Trick formulation
- Different types of Kernel

$$\Rightarrow \begin{aligned} X &= \{u_1, u_2, u_3, \dots, u_n\} \\ Y &= \{y_1, y_2, y_3, \dots, y_n\} \end{aligned} \quad \begin{matrix} \text{input data} \\ \text{or} \\ \text{training data} \end{matrix}$$

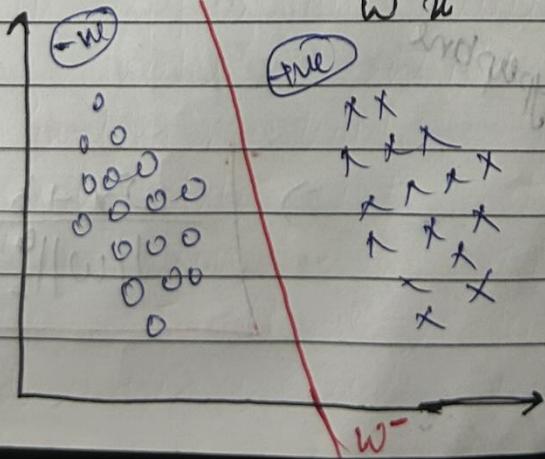
$y_i \in \{-1, 1\}$

Key idea - Separate data with maximum margin hyperplane

$$w^T u + b = 0 \quad (\text{eqn of hyperplane})$$

Hyperplane - having 2 or more features

$$[\theta_0 + \underbrace{\theta_1 u_1^{(1)} + \theta_2 u_2^{(2)} + \theta_3 u_3^{(3)} + \dots }_{w^T u}]$$

 \downarrow b 

{ Hyperplane -
The best decision
boundary }

(1) \rightarrow subscript (sample)

$x_i \rightarrow$ subscript (feature)

Optimal Hyperplane

$$\bullet w^T u + b > 0$$

if $x^{(i)} \in$ true class

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

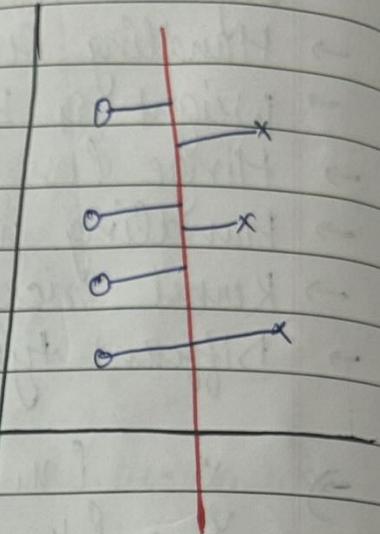
if $x^{(i)} \in$ -ve class

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

Weight matrix

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Input matrix



$$g(x) = +1 \quad \text{if } x \geq 0$$

$$g(x) = -1 \quad \text{if } x < 0$$

Distance

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

hyperplane

$$\Rightarrow \frac{w^T u + b}{\|w\|}$$

$$\boxed{\sqrt{w_1^2 + w_2^2 + \dots + w_n^2}}$$

Euclidean

norm of weight

L2 norm

$$\boxed{\frac{w^T u + b}{\|w\|_2}}$$

Goal - To maximize the min. distance of point from the hyperplane

$$e^{(i)} = w^T u^{(i)} + b$$

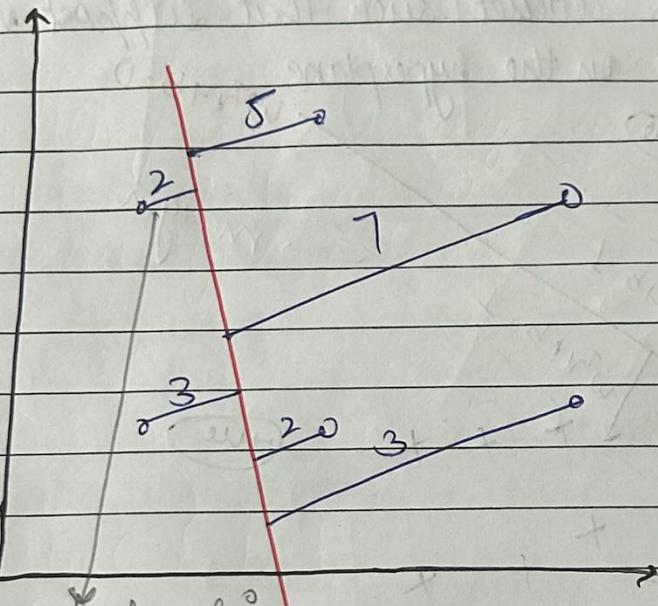
$$\|w\|_2$$

Maximize
min distance

$$g_i = \min_{i=1, \dots, n} e^{(i)}$$

This is how we find
the distance

Now, try to formulate
the goal of SVM using
this notation.



Yeh jo kam hai,
unha distance increase kina hai
takki line distant ho ske & misclassified
points correctly classify ho ske.

Optimisation | Sum Objective

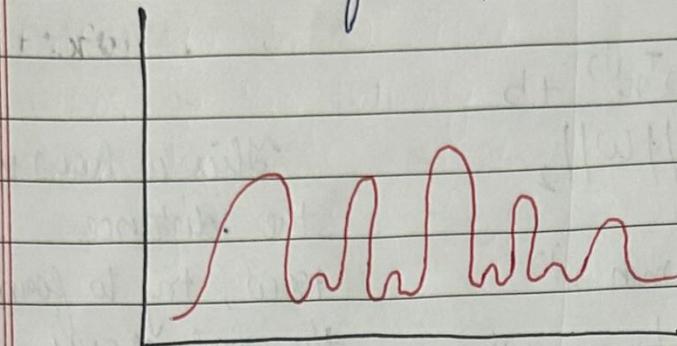
max
 w, b
such that

Absolute
distance

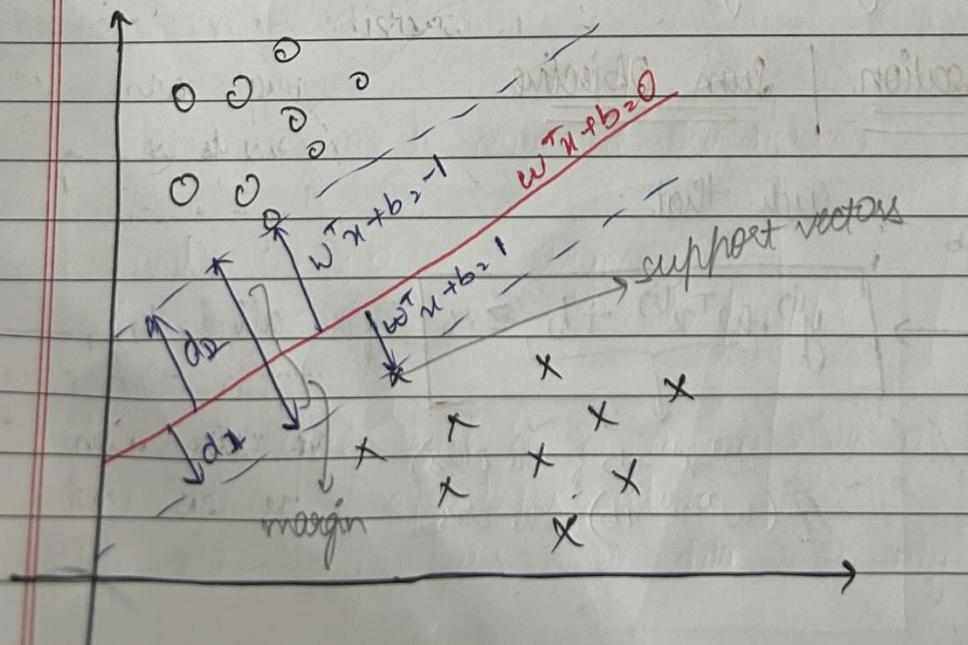
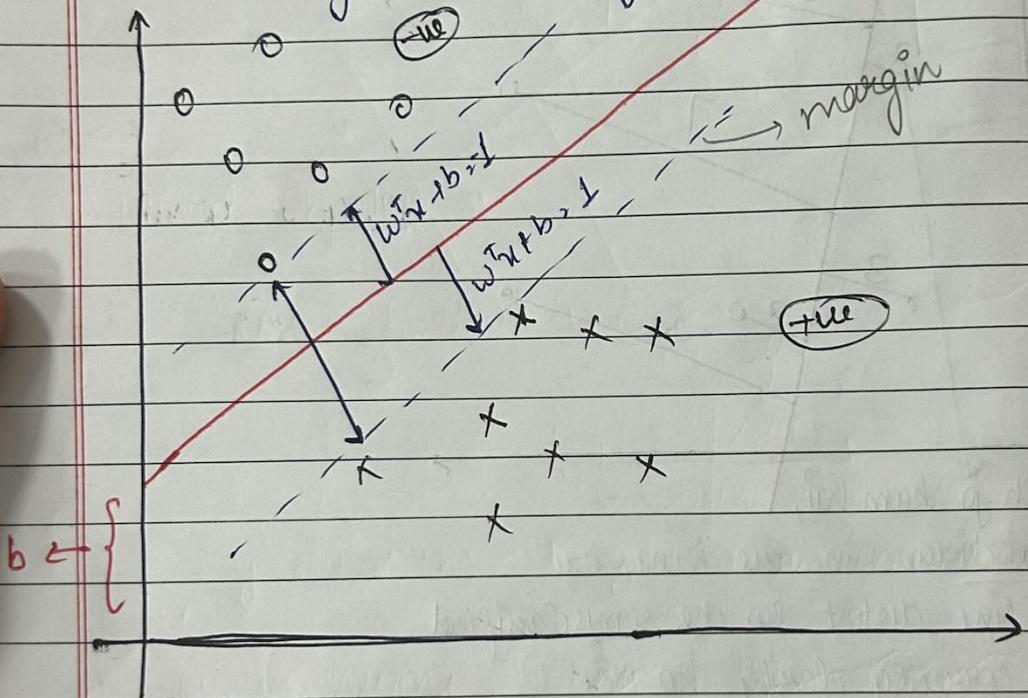
$$y^{(i)} \cdot \underbrace{w^T u^{(i)} + b}_{\downarrow} \geq g \quad \text{fix all } i = 1, 2, \dots, m$$

$w^T u^{(i)} + b$ (true class)
 $-(w^T u^{(i)} + b)$ (wrong class)

* The above written absolute distance ($|y^{(i)} + w^T x^{(i)} + b| \geq \gamma$) is a non-convex function:



→ Renormalise our dataset such that support vectors should always lie on the hyperplane $w^T x + b = 0$



→ Goal - Maximize d such that $y_i(w^T x^{(i)} + b) \geq \epsilon$

$$d_1 = \frac{|w^T x_i + b|}{\|w\|}, \quad \frac{1}{\|w\|} \quad \left\{ \because w^T x_i + b = 1 \right\}$$

$$d_2 = \frac{|w^T x_i + b|}{\|w\|} = \frac{1}{\|w\|}$$

$$d = d_1 + d_2 = \frac{1}{\|w\|} + \frac{1}{\|w\|} = \frac{2}{\|w\|} \rightarrow \text{L-1 norm}$$

Maximise

$$\therefore \text{we minimize } \frac{\|w\|}{2} \quad \left[\begin{array}{l} \text{reciprocal issliye kiya} \\ \text{kyuki hum hamcha minimize} \\ \text{karte aa she hain} \end{array} \right]$$

Final SVM Objective $\rightarrow \min \left\{ \frac{1}{2} \|w\|^2 \right\}$ convex function

$$\text{Final SVM Objective} \rightarrow \min \left\{ \frac{1}{2} \|w\|^2 \right\}$$

$$y_i(w^T x^{(i)} + b) \geq 1$$

This is linear constraint \rightarrow Quadratic solver

Lagrangian duality
(converts constrained
optimisation to unconstrained
optimisation)

min slope = 0

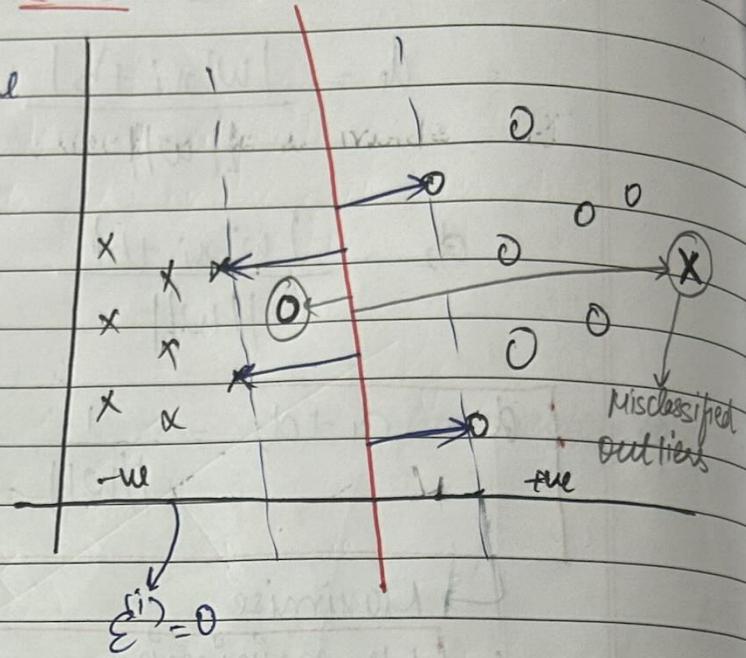
Pegazos - linear constraint
ko kaise unconstrained
dena ekte hain

e.g. $n=2$
(breaks the
cone)

Handling Outliers in SVM

Allow one to do some error on the training examples

$$(x_i, y_i) \rightarrow \varepsilon^{(i)}$$



→ Primal Objective →

$$\text{Loss} = \frac{1}{2} \mathbf{w} \mathbf{w}^T + C \sum_{i=1}^m \zeta^{(i)}$$

$\zeta^{(i)}$ ↓
Zeta
Error

$$y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \zeta^{(i)} \quad \text{Allow some error}$$

If you are having 2 planes so their distance will be,

$$\text{distance} = \frac{2}{\|\mathbf{w}\|}$$

now we have to minimize this distance

$$\text{find } \mathbf{w} \text{ & } b \quad \text{argmin}_{(\mathbf{w}, b)} \frac{2}{\|\mathbf{w}\|}$$

- This distance will be maximum when all points are on one side of $\mathbf{w}^T \mathbf{x} + b$, others at $\mathbf{w}^T \mathbf{x} + b$
- This works only in cases when data is linearly separable

$$\underset{w,b}{\operatorname{argmax}} \frac{2}{\|w\|} \quad \left. \right\} \rightarrow \text{Hard Margin SVM}$$

(does not work on non-linearly separable data)

Date / /
Page No.

② Hard Margin SVM { also called }

sol, Soft Margin SVM \rightarrow for non-linearly separable
(the problem is of outliers)

$$③ \boxed{\underset{w,b}{\operatorname{argmin}} \frac{\|w\|}{2} + C \cdot \frac{1}{n} \sum_{i=1}^n \xi_i^2} \rightarrow \text{hyperparameters}$$

↓
penalty

Soft Margin SVM (bcz we added error)

$$f(w) \rightarrow \max \rightarrow \min \left(\frac{1}{f(w)} \right)$$

Why minimising now?

$$④ \boxed{\frac{1}{n} \sum_{i=1}^n \xi_i} \rightarrow \text{Average distance of misclassified points}$$

(to those points which are not correctly classified, that's distance average it is)

⑤ Pegasos Algorithm for unconstrained optimisation

$$\begin{aligned} w^T u_1 + b &= 1 \\ w^T u_2 + b &= -1 \end{aligned}$$

$$\frac{w^T (u_1 - u_2)}{\|w\|} = \frac{2}{\|w\|}$$

* SVM - find w, b such that where $\frac{\|w\|}{2}$ is minimum

$$\therefore \boxed{(u_1 - u_2) = \frac{2}{\|w\|}}$$

16/2/23

Date / /
Page No.

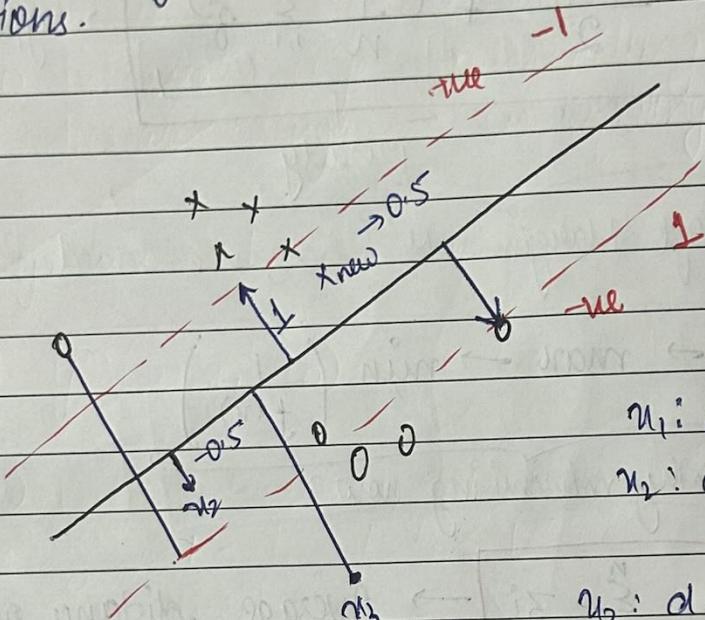
Pegasos Algorithm for Unconstrained Optimisation

→ Objectives-

- i) Develop a gradient descent algorithm for SVM
- ii) Introduce sub-gradients for convex but non-differentiable cost functions.

$y_i w^T x > 0$
then correctly
classified

if
 $y_i w^T x > 1$
misclassified



$$\begin{aligned} u_1: d &> 0.5 (1 - 0.5) \\ u_2: d &> -0.5 \\ 1 - (-0.5) &= 1.5 \\ u_3: d &> -1.5 \\ 1 - (-1.5) &= 2.5 \end{aligned}$$

$$\underset{w,b}{\operatorname{argmin}} \frac{\|w\|}{2}$$

distance
from
correct
plane
for
margin

* SVM mein non-convex function hota hai, we make hyperplanes

→ Non Convex

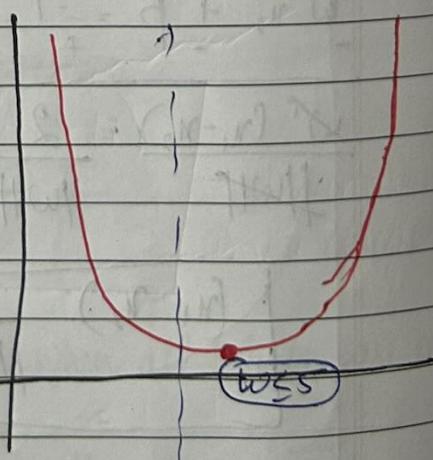


Convex with Constraints



Convex with Constraints + Outliers

$$\left[\text{Min } \frac{\|w\|}{2} \right] \text{s.t. } y^i (w^T x_i + b) \geq 1$$



⇒ Due to constraints, we're unable to differentiate the function

* Constraint function



conversion

Unconstrained function

⇒ When it will get converted, the differentiation problem will get solved. Gradient Descent will be able to optimise

→ When $d > 1$, this is the case of misclassification
(This is the error)

→ d should be less than 1 (for correct classification)

⇒ Misclassified distance from correct plane:

1- (distance of misclassified points from correct plane)

$\zeta^{(i)}$
(zeta)

Note: ζ_i will increase as point is further away from correct plane

$\zeta_i \rightarrow$ is for outliers (~~error~~)

$y_i w^T x_i > 0 \Rightarrow$ correctly classified

Class label. distance ≥ 0

* We have to tune the hyperparameter.

Date	1/1/
Page No.	

→ Loss Function / Objective

$$w, b = \underset{\text{Regulariser}}{\operatorname{argmin}} \frac{\|w\|}{2} + C \cdot \underset{\text{Hyperparameter/penalty}}{\underset{\text{loss function}}{\sum_{i=1}^n z_i}}$$

- * If C too increase Krenge \rightarrow overfitting
If C too decrease Krenge \rightarrow underfitting

$$\Rightarrow d = 1 - y_i w^\top x_i \geq 1$$

$$1) z^{(i)} = 1 - y^{(i)}(w^\top x_i + b)$$

$$2) z^{(i)} \geq 0$$

$z^{(i)} \geq 1 - t_i \text{ if } t_i \leq 1$

$$z^{(i)} \geq 0$$

if $t_i \geq 1$

$$2) z^{(i)} = 1 - t_i \text{ if } t_i \leq 1$$
$$z^{(i)} = 0 \text{ if } t_i > 1$$

$$z^{(i)} = \max(0, 1 - t_i)$$

Loss Function

$$\min_{w,b} \frac{1}{2} \frac{\|w\|^2}{2} + C \sum_{i=1}^m \max(0, 1 - t_i)$$

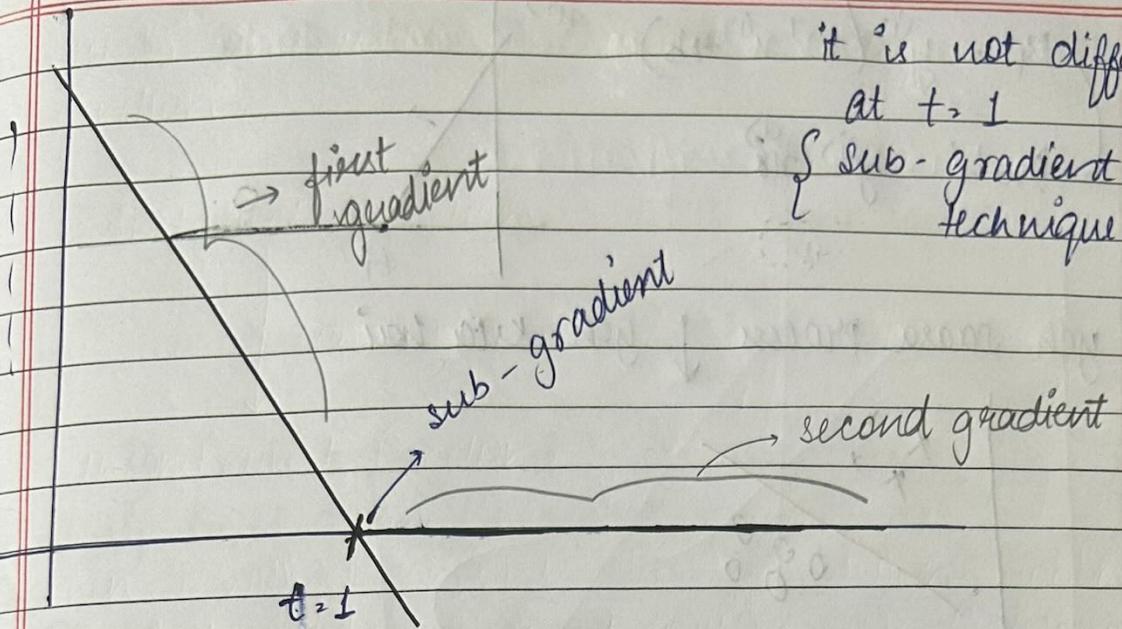
$$\text{where, } t_i = y^{(i)}(w^\top x^{(i)} + b)$$

→ unconstrained convex optimisation

↳ Hinge loss

it is not differentiable
at $t \rightarrow 1$

{ sub-gradient
technique }



$$w + c \sum_{i=1}^m \nabla_w [\max(0, 1-t_i)]$$

$$= w + c \sum_{i=1}^m \nabla_w [\max(0, 1-t_i)]$$

$$\nabla_w t_i \leftarrow \begin{cases} 1 & \text{if } t_i \geq 1 \\ -1 & \text{if } t_i < 1 \end{cases}$$

CHAIN RULE

$$y = w + c \sum_{i=1}^m \begin{cases} 0 & \text{if } t_i \geq 1 \\ -1 & \text{if } t_i < 1 \end{cases} g^{(i)} x^{(i)}$$

Differentiation of Hinge Loss

$$\Delta w_k = w + c \cdot \sum_{i=1}^m \nabla_w [\max(0, 1-t_i)]$$

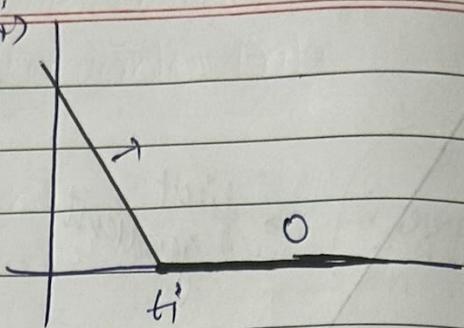
$$= w + c \sum_{i=1}^m \nabla_t [\max(0, 1-t_i)] \nabla_w t_i$$

CHAIN RULE

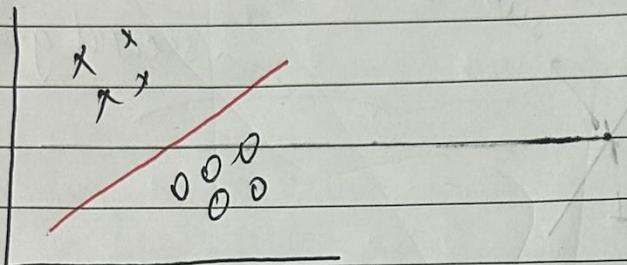
$$= w + c \cdot \sum_{i=1}^m \begin{cases} 0 & \text{if } t_i \geq 1 \\ -1 & \text{if } t_i < 1 \end{cases}$$

$$t_i = y^{(i)}(w^T x^{(i)} + b)$$

$$\Delta w t_i \rightarrow y^{(i)} x^{(i)}$$



yeh saara process \downarrow yeh kita hai



~~20/2/23.~~

SVM: Weight & Bias Update Rule {The final update rule}

$$L(w) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max(0, 1 - t_i)$$

$$\text{where, } t_i = y^{(i)}(w^T x^{(i)} + b)$$

Convex loss function will be Gradient Descent to find optimal values of w & b .

$$i) w = w - n \cdot \Delta_w L$$

$$ii) b = b - n \cdot \Delta_b L$$

derivative of Loss w.r.t w

[(i) & (ii) are two equations we have to find]

$$\Rightarrow \Delta_w L = w + C \begin{bmatrix} 0 & t_i \geq 1 \\ -1 & t_i < 1 \end{bmatrix} y^{(i)} x^{(i)}$$

(After deriving eq (i) & (ii) w.r.t w & b resp.)

$$\rightarrow D_{bt} = 0 + c \begin{bmatrix} 0 & t_i \geq 1 \\ -1 & t_i < 1 \end{bmatrix} y_i$$

(After deriving eq (i) & (ii) wrt w & b resp.)

$$D_{wt} = y_i (w^T x^{(i)} + b)$$

- i) $D_{wt} \rightarrow y^{(i)} w^{(i)}$ in wrt w
- ii) $D_{bt} \rightarrow y^{(i)}$ wrt b

Final weight update rule

$$w = w - n \left[w + c \begin{bmatrix} 0 & t_i \geq 1 \\ -1 & t_i < 1 \end{bmatrix} y_i x_i \right]$$

expanding above eqn
Simplifying:

$$\rightarrow w = w - nw \begin{bmatrix} 0 & t_i \geq 1 \\ nc[y_i x_i] & t_i < 1 \end{bmatrix}$$

$$\rightarrow b = b - n \left[c \begin{bmatrix} 0 & t_i \geq 1 \\ -1 & t_i < 1 \end{bmatrix} y_i \right]$$

$$\begin{cases} b = b + 0 & \text{if } t_i \geq 1 \\ b = b + ny_i & \text{if } t_i < 1 \end{cases}$$

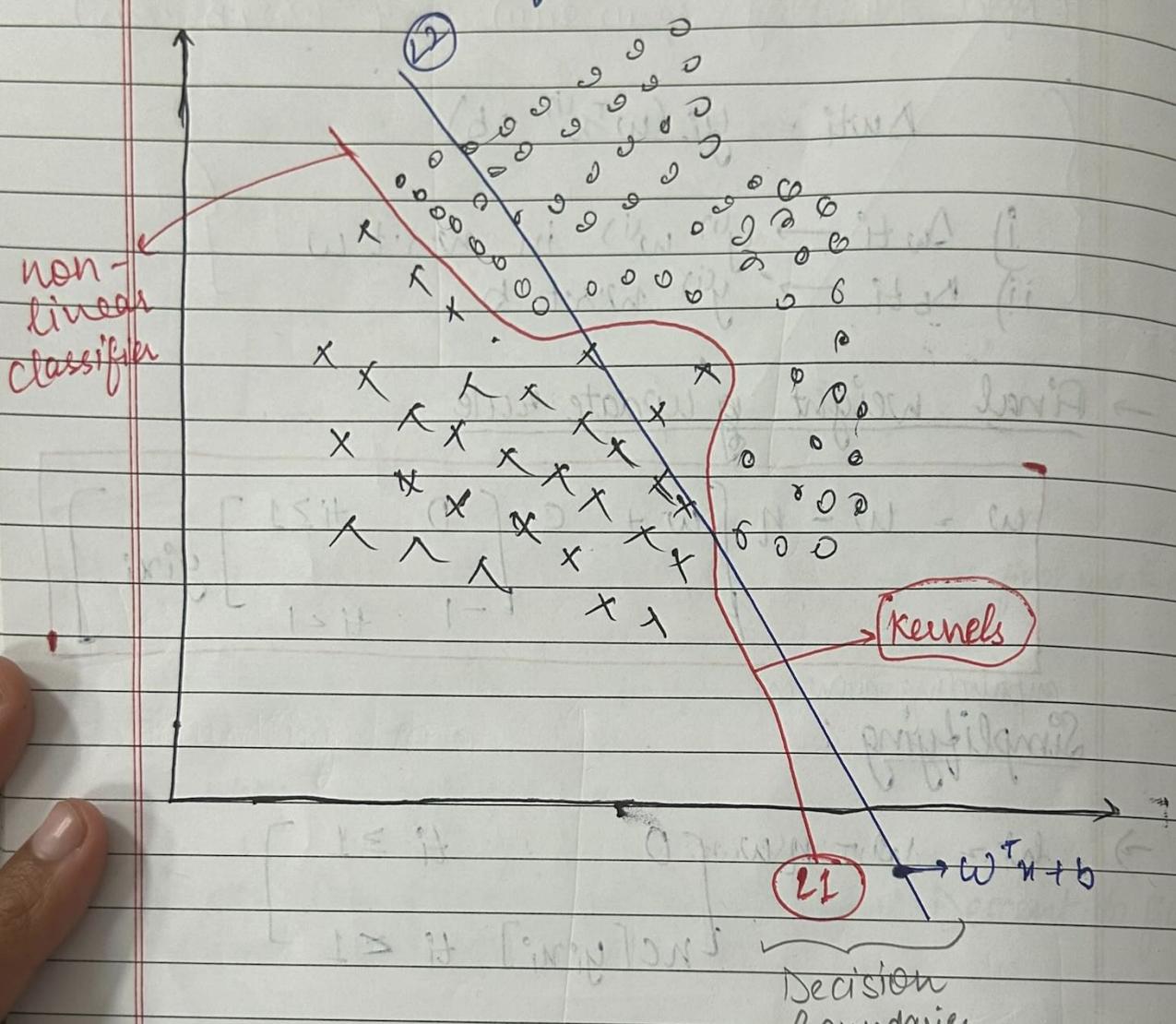
* Margin maximisation works well on linearly separable data.

SVM is for both regression & classification

Date / /
Page No.

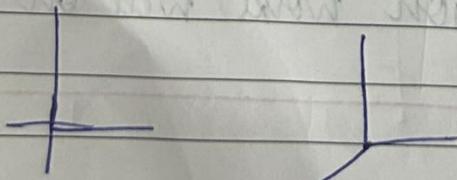
Handling Non-linear Separable Data

Non-Linear Classification

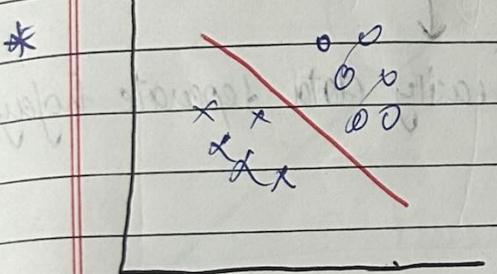
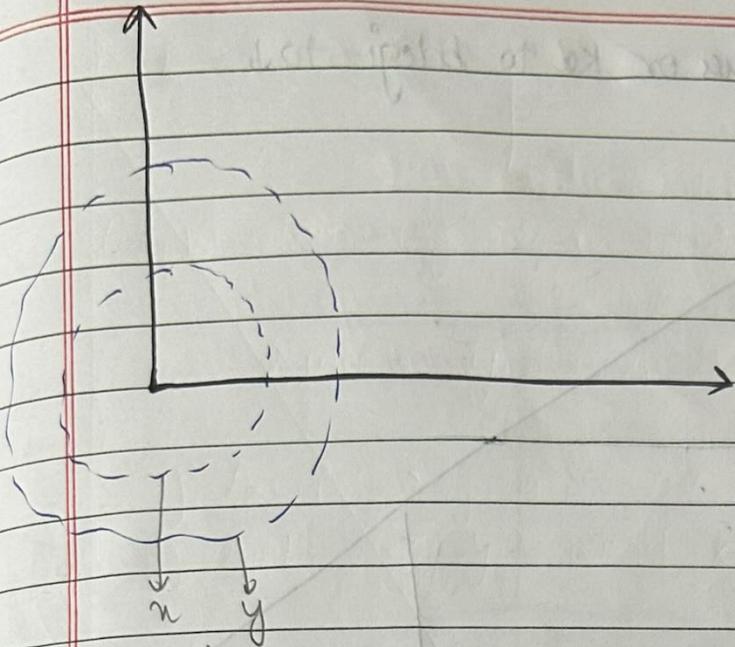


- Send this data to higher dimension & separate the classes using linear classifier.
- Key Idea - Project Data to High Dimension

$$x_2 [x_1, x_2] \xrightarrow{\text{Transform}} x_2 [x_1, x_2, x_3]$$

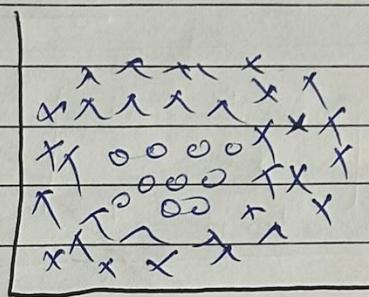


$$x_4 = x^2$$



Linear separable
data

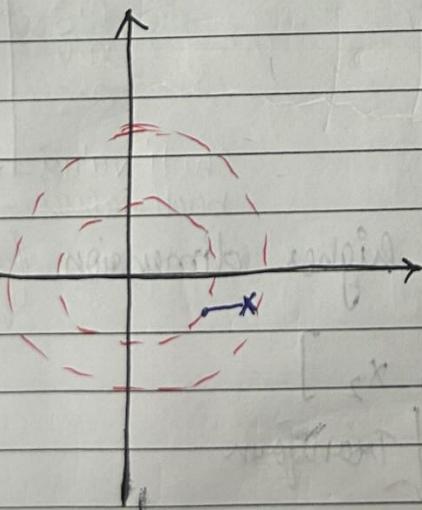
(Margin maximisation
is used)



Non-linear separable
data

(Kernel trick is
used)

* Kernel trick uses math to separate data



~~Original~~ \rightarrow y_1

0.1

0.2

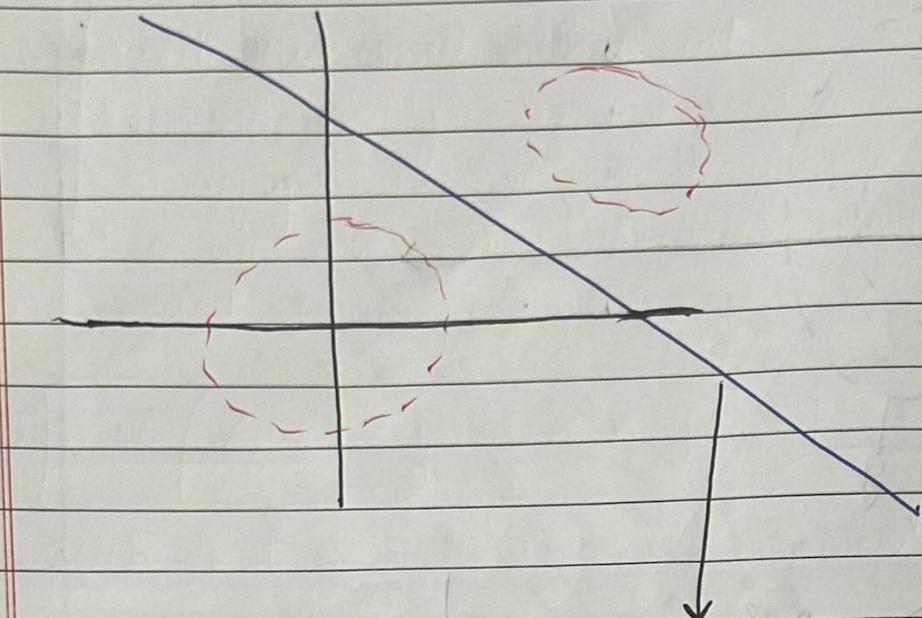
y (y_2)

2

4

$$\text{abst } y_3 = y_1^2 + y_2^2$$

x_1^2 kرنge to values or ka to bdegi to j



ab easily data separate hoga yeha

~~22/2/23~~

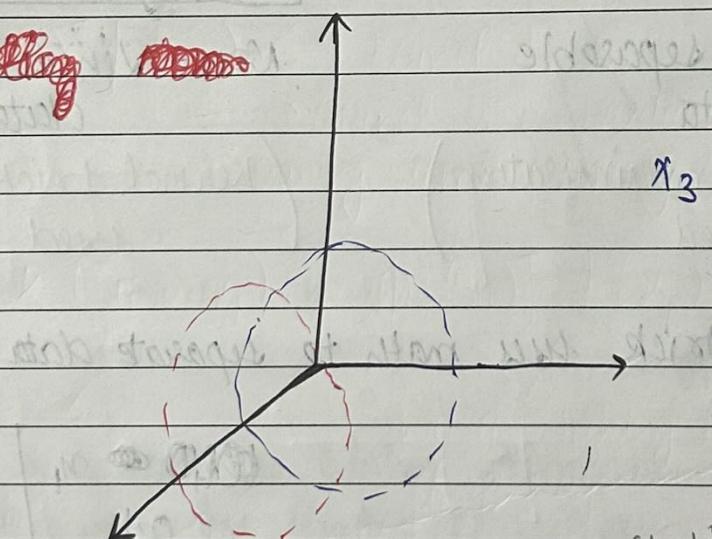


Projecting

points

separating now
will

$$x_3 = x_1^2 + x_2^2$$



w_1	w_2
0.2	2
0.3	4
0.4	6
0.5	8

choti values → all choti
badi values → all badi

Key Idea - Project data to higher dimension (after sq.)

$$x = [x_1, x_2]$$

↓ Transform

$$x = [x_1, x_2, x_3]$$

$$\Rightarrow x_3 = x_1^2 + x_2^2$$

DR

(n.)

→ (n₂): Back points are
back · bags in plane
in upper bags

→ Plane

→ chose points our
chose hogye

Made easy to
separate

'Kernel Trick' based formulation

Ques) How can you apply sum to non linearly separable data?

→ The first change you have to make to the formulation is

$$\min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n y_i (w^T x_i + b) \geq 1 - y_i$$

Loss function

$$u_i \rightarrow \phi(u_i)$$

will transform x_i & change x_i

→ $\phi(x_i)$

→ but this comes with a cost
 → computationally expensive because you are going to apply transform over every data point
 → so we have often formulation which is SVM based on Lagrangian

→ $\max \left(\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \phi(x_i)^T \phi(x_j) \right)$

such that

$$\sum \alpha_i y_i = 0 \quad (\text{for all samples})$$

This is the result

→ $\phi(x_i)^T \phi(x_j) \rightarrow \text{Kernel trick}$

→ kernel is a particular function having some property.

→ $k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$

We are taking product of 2 projected vectors

→ RBF

Polynomial
Sigmoid
Linear

} Different types of Kernel

Different kernels decision boundary bivariate Lai

Date / /
Page No.

→ RBF :- Radial Basis Kernel

$$K(x_i, x_j) = e^{-\gamma \|x_i - x_j\|^2} \quad \{\gamma \text{, gamma}\}$$

↓
This does transformation

coefficient γ
(only control kya ta hoi)

→ Polynomial Kernel :-

$$K(x_i, x_j) = (\gamma x_i^T \cdot x_j + r)^c$$

degree of
polynomial

→ Sigmoid -

$$K(x_i, x_j)$$

$$K \text{ of } x_i, x_j = \frac{1 - e^{-2(y x_i \cdot x_j + r)}}{1 + e^{-2(y x_i \cdot x_j + r)}}$$

Ques) Why these kernels helps in simple calculation?

→ ϕ is a map from n-dimension to m-dimension space

Compute $\phi(u), \phi(y)$,

do dot product

→ this will be expensive, isliye Kernel trick dekha

→ Without kernel trick,

Eg) suppose $x = (1, 2, 3)$
 $y = (4, 5, 6)$

$$f(x) = (1, 2, 3, 4, 6, 3, 6, 9)$$

$$f(y) = (16, 20, 24, 20, 25, 30, 30, 36)$$

$$= f(x), f(y) \geq 16 + 40 + 72 + 40 + 100 + 180 + 72 + 180 + \\ 324 = 1024$$

↓
lot of algebra, mainly because
it is mapping from 3D to 4D

⇒ With Kernel

$$K(x, y) = (4 + 10 + 18)^2 = 32^2 = 1024,$$

same result, but calculation is easier

(took dot product of x, y)

$$(1, 2, 3) \quad (4, 5, 6)$$

$$(1 \times 4) + (2 \times 5) + (3 \times 6)$$

$$= (4 + 10 + 18)^2 = 32^2 = 1024$$

* Suppose we have imbalanced samples

$$n=10$$

Male Female

So we give more weightage to less samples

(Eg - female ko 4 weightage de diya)

↳ Isse overfitting nahi hoti