

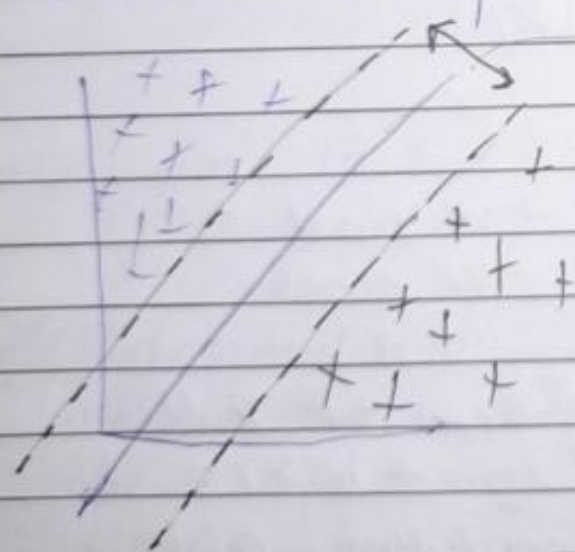
# # Support Vector Machine

- It is a classifier that works both on linearly and non-linearly separable data.



Finds an optimal hyperplane, that best separates our data so that the distance from nearest points in space to itself (also called margin) is maximized.

These nearest points are called Support Vectors.



goal is to maximize this distance

$$ax + by + c = 0$$

dist  $(x, y)$

$$ax + by + c = \text{dist} \times \sqrt{a^2 + b^2}$$

A hyperplane is a plane of  $n-1$  dimension in  $n$  dimensional feature space, that separates the two classes.  $2d \rightarrow$  line  $3d \rightarrow$  plane

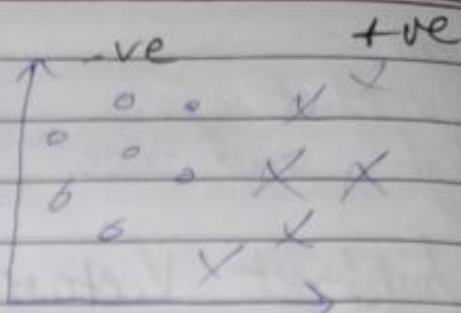
**Max Margin hyperplane**:- The optimal hyperplane that separates our data so that the distance / margin from nearest point (called support vectors) in space to itself is maximized.

## # SVM - Mathematically

$$X = \{x_1, x_2, \dots, x_m\}$$

$$Y = \{y_1, y_2, \dots, y_m\}$$

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



Key idea = separate data with binary classification  
maximum margin

Hyperplane

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

vector

bias intercept

$$a_0 + a_1 x_1 + a_2 x_2 + a_3 x_3 + \dots$$

$\downarrow$

$b$

$w^T x$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_n \end{bmatrix}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$

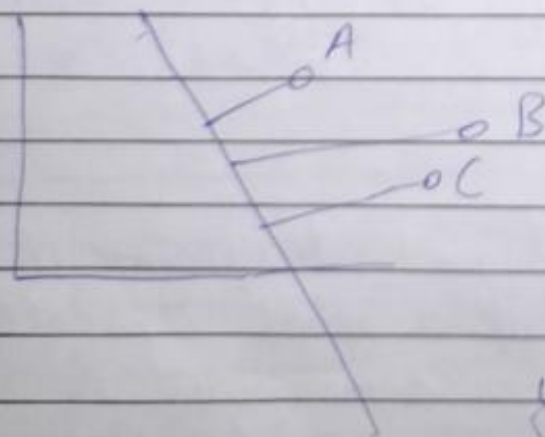
Optimal Hyperplane

$$w^T x + b \geq 0$$

if  $x' \in +ve$  class

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

if  $x' \in -ve$  class



$$d_A = w^T x(A) + b > 0$$

$$d_B = w^T x(B) + b > 0$$

more confident abt prediction at B

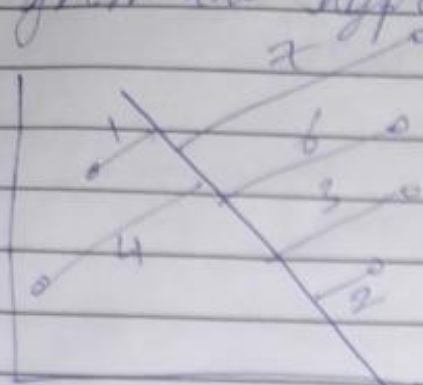
$$y_{pred} = g(w^T x + b)$$

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

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



goal :- To maximize the minimum distance of points from the hyperplane



$$r^i = \frac{w^T x^i + b}{\|w\|}$$

min dis

$$r = \min_{i=1 \dots m} (r^i)$$

All the points should have at least  $r$  distance

$$\left. \begin{array}{l} w^T x + b > 0 \quad \text{if } x +ve \\ w^T x + b < 0 \quad \text{if } x -ve \end{array} \right\} \text{Non convex}$$

Optimization / SVM objective :-

absolute dist of any pt from hyperplane

$$\max_{r, w, b}$$

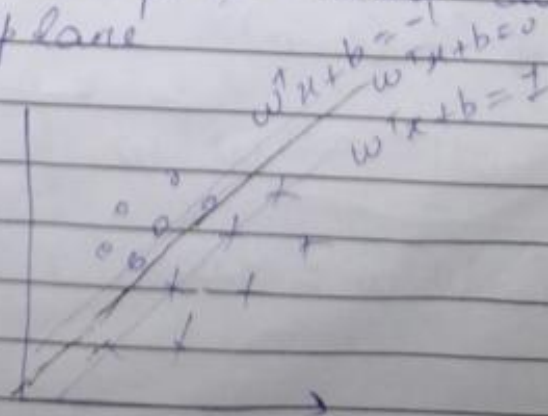
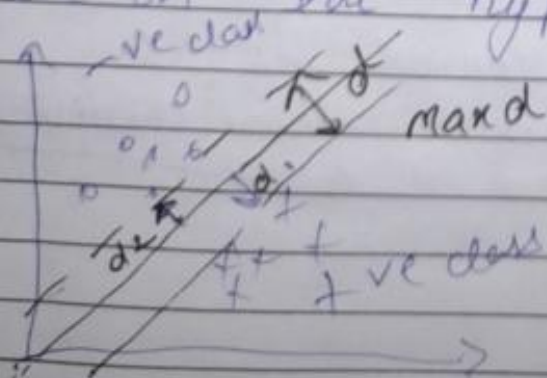
$$\text{such that } y^i (w^T x^i + b) \geq r \text{ for all } i=1 \dots m$$

class

$$\max(r) = \min_{i=1 \dots m} \left( \frac{w^T x^i + b}{\|w\|} \right)$$

Reformulation  $\Rightarrow$

We want to normalise our dataset such that the support vectors (nearest pts) should always lie on the hyperplane



Support vectors should be equidistant  
goal is to maximise distance such that  
all pts are correctly classified

goal:  $y_i (w^T x^i + b) \geq 2$

$$d_1 = \frac{|w^T x_p + b|}{\|w\|} = \frac{1}{\|w\|}$$

$$d_2 = \frac{|w^T x_{-ve} + b|}{\|w\|} = \frac{1}{\|w\|}$$

$$d = d_1 + d_2 = \frac{2}{\|w\|} \rightarrow \text{maximize distance}$$

SVM objective minimize  $\frac{1}{2} \|w\|^2$  inversion

min  $\left\{ \frac{1}{2} \|w\|^2 \right\}$

$$y_i \frac{(w^T x^i + b)}{\|w\|} \geq \frac{1}{\|w\|} \rightarrow \text{normalization term}$$

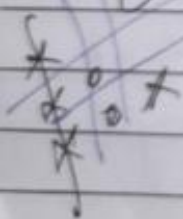
minimize  $\frac{1}{2} \|w\|_2^2 = \sqrt{\sum w_i^2}$

st  $y_i (w^T x^i + b) \geq 1$   
 $\forall i \in (1..n)$

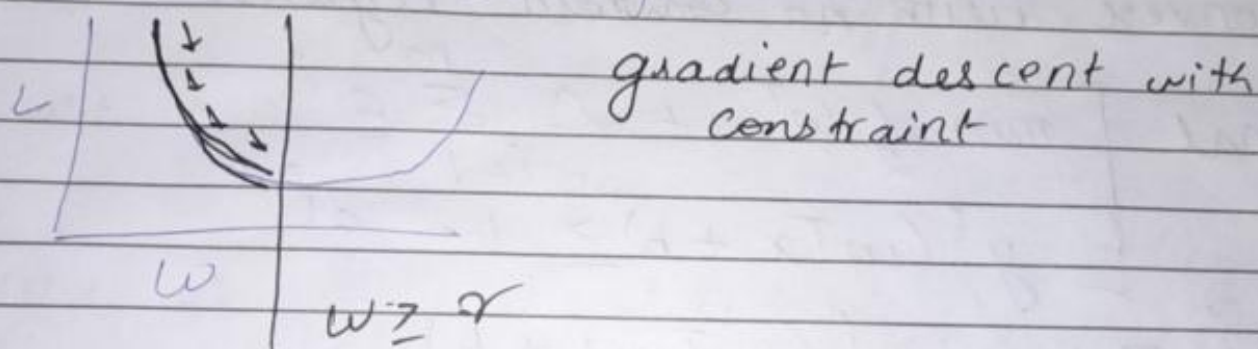
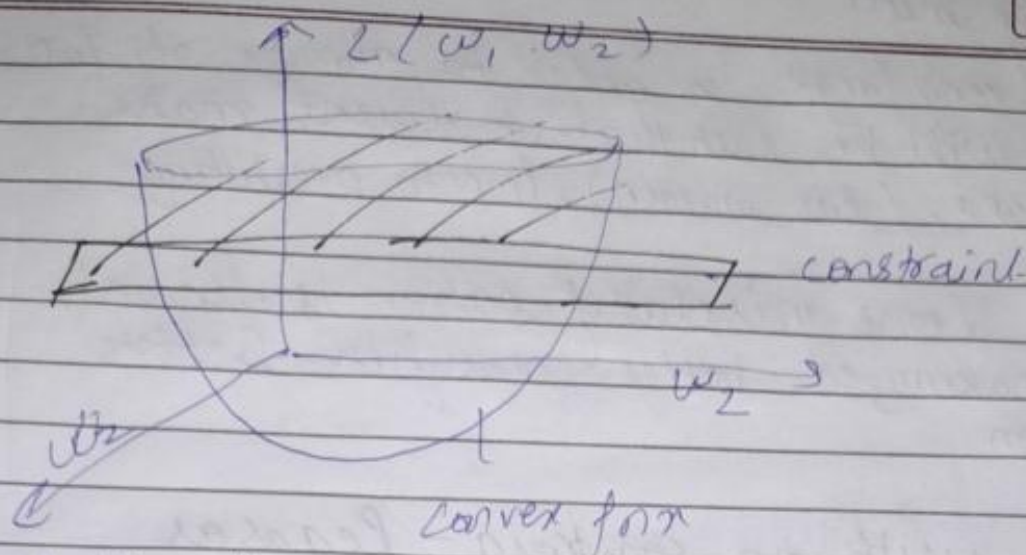
$$\|w\|_2^2 = w^T w$$

$$\nabla_w w^T w = w$$

convert fcn  
with linear constraint



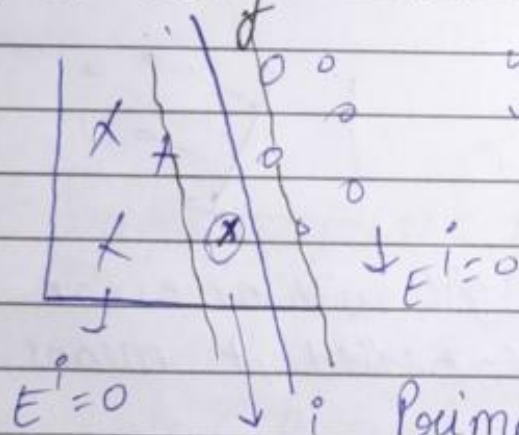




## # Handling Outliers in SVM

allow our algorithm (objective) to do some misclassification on some training example.

$(x_i, y_i) \rightarrow E^i$   
add penalty to SVM objective.



Primal Objective :-

convex with  
constraint +  
outlier

$$\frac{1}{2} w \cdot w^T = \text{loss} + c \sum_{i=1}^m E^i$$

$$y^i (w^T x^i + b) \geq 1 - E^i$$

New formulation

$$\frac{1}{2} w \cdot w^T + c \sum_{i=1}^m E_i \quad [E^i = 0]$$

$$y^i (w^T x^i + b) \geq 1 - E^i$$

$C \rightarrow$  hyperparameter

Date :

If  $C \rightarrow \infty$  / very large in order to minimize obj for optimization will be such that it doesn't make any mistake. (big margin). Prone to overfitting

$C \rightarrow$  small some misclassification is allowed, margin is maximized. Better classifier. Better generalization

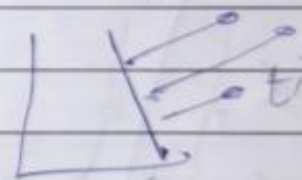
# Convex with no constrain - Pegasas

Primal 
$$\begin{cases} \min \frac{1}{2} \omega^T \omega + C \sum_{i=1}^m \xi_i \\ y_i (\omega^T x_i + b) \geq 1 - \xi_i \end{cases}$$

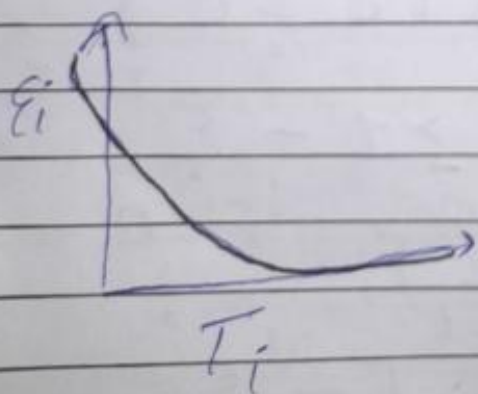
①  $\xi_i \geq 1 - y_i (\omega^T x_i + b)$

②  $\xi_i \geq 0$

$\xi_i \geq 1 - t_i$        $\xi_i \geq 0$



if  $t_i \geq 1$  pt is lying on right with no error  
no since  $\xi_i \geq 1 - t_i$  is not satisfied it means  
it has no error



$\xi_i = 1 - t_i$  if  $T_i \leq 1$

$\xi_i = 0$  if  $t_i \geq 1$

$\xi_i = \max(0, 1 - t_i)$



Loss / objective  $\rightarrow$ 

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

where  $t_i = y_i (w^T x_i + b)$ 

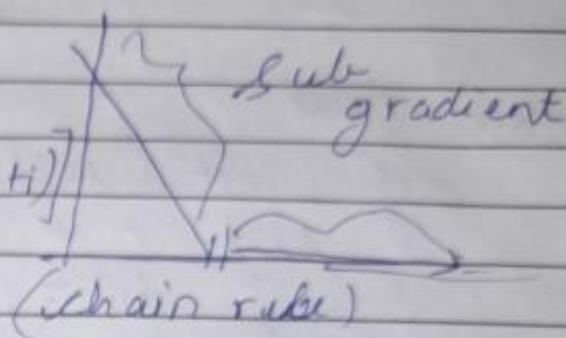
# Applying gradient descent

 $w = \text{init}()$ 

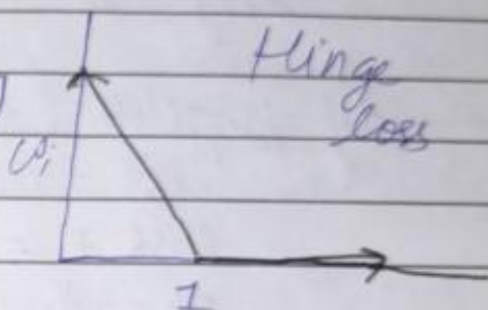
$$w = w - \eta \nabla_w L$$

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

$$= w + c \sum_{i=1}^m \frac{\partial y}{\partial t} \cdot \frac{\partial t}{\partial w}$$



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



$$t_i = y_i (w^T x_i + b)$$

$$\nabla t_i = y_i x_i$$

# weight update rule

$$w = w - \eta \nabla_w L$$

$$b = b - \eta \nabla_b L$$

$$\nabla_w L = w + c \begin{bmatrix} 0 & t_i \geq 1 \\ -1 & t_i < 1 \end{bmatrix} y_i x_i$$

$$\nabla_b L = 0 + c \begin{bmatrix} 0 & t_i \geq 1 \\ -1 & t_i < 1 \end{bmatrix} \nabla_b t_i$$

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

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

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

class SVM:

def \_\_init\_\_(self, c=1.0):

self.c = c

self.w = 0

self.b = 0

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

def hinge\_loss(self, w, b, x, y):

loss = 0.0

loss += 0.5 \* np.dot(w, w.T)

m = x.shape[0]

for i in range(m):

t\_i = y[i] \* np.dot(w, x[i].T) + b

loss += self.c \* max(0, 1 - t\_i)

return loss[0][0]

$$w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

def fit(self, x, y, batch\_size=100, lr=0.001):

no\_of\_feature = x.shape[1]

no\_of\_sample = x.shape[0]

η = learning\_rate

c = self.c

# Init the model parameters

w = np.zeros((1, no\_of\_feature))

bias = 0

print(self.hinge\_loss(w, bias, x, y))



# weight and bias update rule

losses = 0

for i in range(max\_itr):

l = self.hingeloss(w, bias, X, y)

# Batch Gradient Descent with random shuffle  
for batch\_start in range(0, no\_of\_samples, batch\_size):

# assume 0 gradient for the batch

gradw = 0

gradb = 0

# iterate over all cx in mini batch

for j in range(batch\_start, batch\_start + batch\_size):

if j < no\_of\_samples:

i = ids[j]

t\_i = y[i] \* np.dot(w, X[i].T) +

if t\_i > 1: bias)

gradient < gradw += 0  
gradb += 0  
else

< gradw += c \* y[i] \* X[i]  
gradb += c \* y[i]

# updation

W = W - lr \* W + lr \* gradw

b = b + lr \* gradb

return W, bias, losses

## # SVM Kernel

- apply SVM to non-linearly separable data

For linearly separable

$$\min \int \frac{1}{2} w^T w + c \sum_i \epsilon_i$$

$$y_i (w^T \phi_i + b) \geq 1 - \epsilon_i$$

$$x_i \longrightarrow \phi(x_i)$$

computationally expensive

Another formulation:- Based on Lagrangian

$$\max \left( \sum_i \alpha_i - \frac{1}{2} \sum_i \alpha_i \alpha_j y_i y_j \underbrace{\phi(x_i)^T \phi(x_j)}_{\text{computation}} \right)$$

$$\sum_i \alpha_i y_i = 0$$

Kernel Trick :-

A Kernel is a fn  $k$  which has following prop

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

↑  
projected vector of  $i^{\text{th}}$  vector

- ① RBF Kernel
- ② Polynomial Kernel
- ③ Sigmoid Kernel
- ④ Linear Kernel



① RBF Kernel - Radial Basis Kernel

$$K(x_i, x_j) = e^{-\gamma \|x_i - x_j\|^2} \rightarrow x_i - x_j \cdot (x_i - x_j)^T$$

$\gamma$  = amplitude of freq

② Polynomial Kernel

$$K(x_i, x_j) = (\gamma x_i^T \cdot x_j + c) \rightarrow \text{degree of polynomial}$$

③ Sigmoid Kernel

$$K(x_i, x_j) = \frac{1 - e^{-2(\gamma x_i^T x_j + c)}}{1 + e^{-2(\gamma x_i^T x_j + c)}}$$

Why Kernel?

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

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

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

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

So long!

using Kernel =  $K(x, y) = [(1, 2, 3) \cdot (4, 5, 6)]^2$

$$= [4 + 10 + 18]^2$$

$$K(x, x) = [\langle x, y \rangle]^2 = 1024$$

done so simple