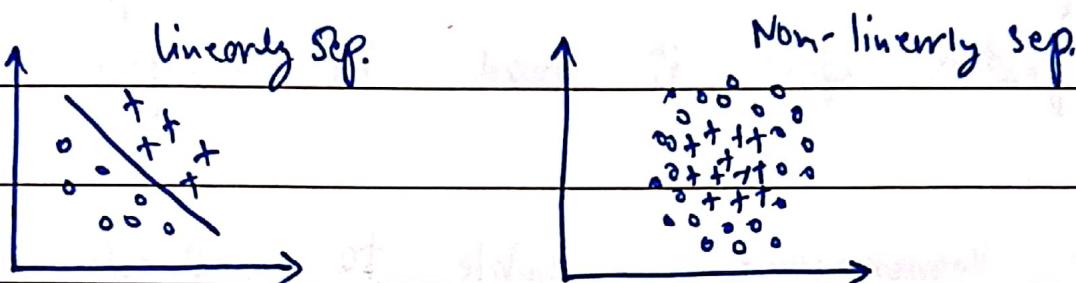
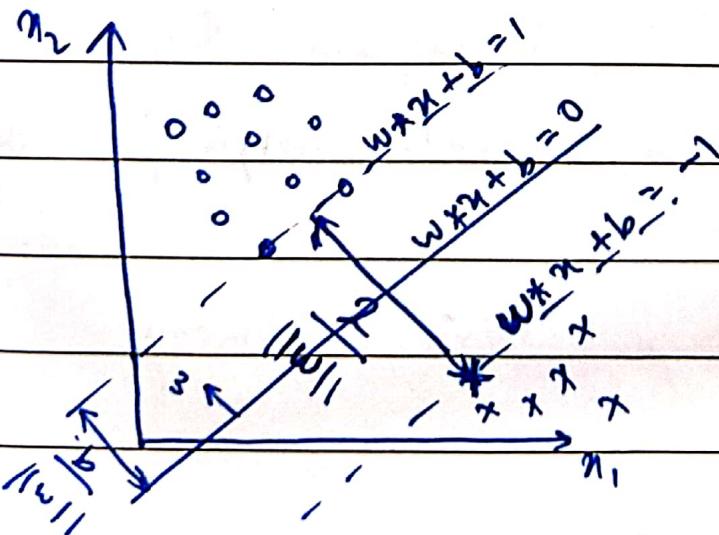


Support Vector Machines (SVM)

- * SVM is a powerful classifier that works both on linearly & 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.



- * for non-linearly^{sep.} data, it uses 'Kernel Trick'

What is Hyperplane?

A hyperplane is plane of $n-1$ dimensions in n dimensional feature space, that separates the two classes. For a 2-D feature space, it would be a line & for a 3-D feature space, it would be a plane and so on.

A hyperplane is able to separate classes if for all points:

$$w_n + b > 0 \quad (\text{for data points in class 1})$$

$$w_n + b < 0 \quad (\text{for data points in class 0})$$

(Note: $w_n + b$ is distance of a pt. from plane)

Goal of SVM: Find hyperplane that separates data with max. margin from closest points.

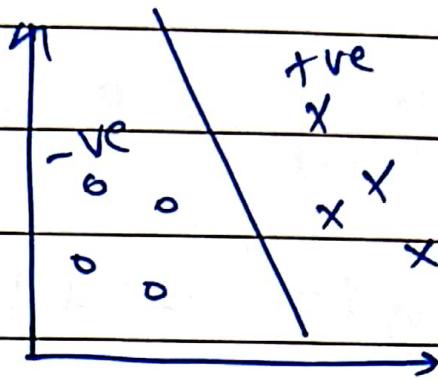
* Example: Say you have a binary classification.

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

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

Each x has n no. of features.

$$y_{(i)} \in \{-1, 1\}$$



Key Idea: Separate data with max. margin.

Hyperplane : $w^T n + b = 0$

$$w^T n = \theta_0 + \underbrace{\theta_1 n^{(1)} + \theta_2 n^{(2)} + \theta_3 n^{(3)} + \dots}_{w^T n} + b$$

$$w = [w_1, \dots, w_n], \quad n = [n_1, \dots, n_n]$$

Optimal Hyperplane :

* $w^T n^{(i)} + b > 0$ if $n^{(i)} \in$ +ve class

$w^T n^{(i)} + b < 0$ if $n^{(i)} \in$ -ve class

* Predictions should be wrong & confident.

More distance of point from hyperplane

means more confident wry if we change

hyperplane then large distance might not change the class.

$$y_{\text{pred}} = g(w^T n + b)$$

$$g(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

* Unlike logistic, we ~~will~~ will not compute probabilities.

* How to find distance

$$w^T n + b$$

$$\bullet (n)$$

$$\text{dis} = \frac{w^T n + b}{\sqrt{w_1^2 + w_2^2 + \dots}}$$

$$= \frac{w^T n + b}{\|w\|_2}$$

(L₂ norm)

L₂ norm is Euclidean distance of point from origin.

Goal : Maximise the minimum distance of point from the hyperplane.

$$\text{Distance, } \gamma^{(i)} = \frac{\mathbf{w}^T \mathbf{x}^{(i)} + b}{\|\mathbf{w}\|_2}$$

$$\text{Min distance, } \gamma = \min_{i=1 \dots m} \gamma^{(i)}$$

So, all the points should have atleast γ distance.

optimization / sum objective :

$$\begin{aligned} & \text{find max value of } \gamma \quad [\max_{\gamma, \mathbf{w}, b}] \\ & \text{such that } \underbrace{y^{(i)} + (\mathbf{w}^T \mathbf{x}^{(i)} + b)}_{i=1 \dots m.} \geq \gamma \quad \text{for all} \\ & \qquad \qquad \qquad \text{(This just gives absolute distance.)} \\ & \qquad \qquad \qquad -1 \neq d = +ve \end{aligned}$$

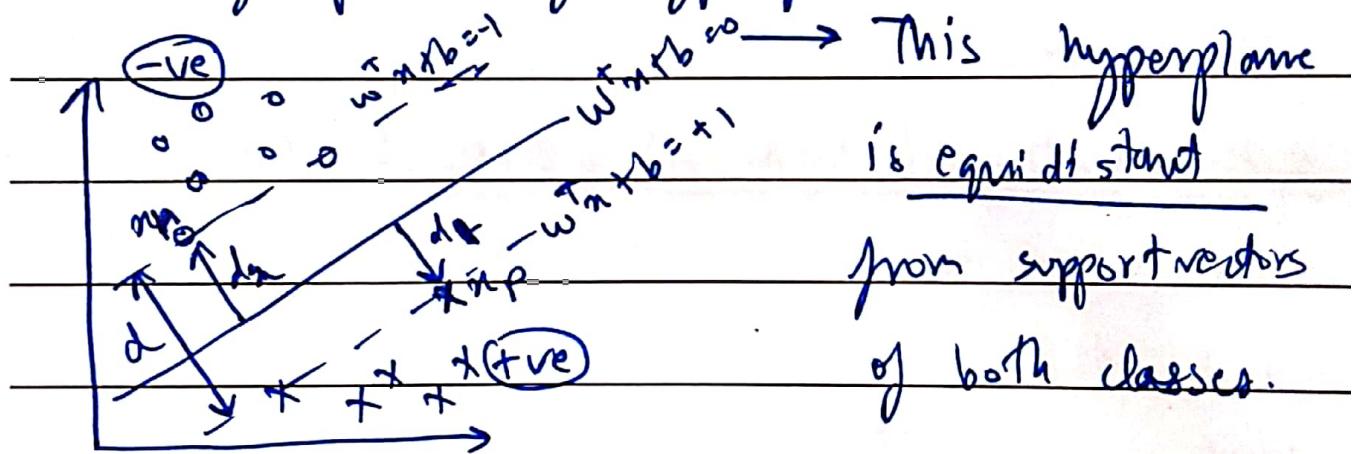
→ This is very hard to solve (non-convex func)
∴ Hence we will convert it to a simple solvable way.

function has many local minima/ maxima, hence cannot find global answer.

Objective or Constrained :

Reformation : We re-normalise our data set such that the support vectors (points closest to the hyperplane) should always lie on the hyper plane.

Using support vectors we maximize the margin of the classifier. Deleting Support vectors will change position of hyperplane.



Goal : Maximize d such that $y_i \cdot w^T n^{(i)} + b \geq \gamma$

$$d_1 = \frac{|w^T n_p + b|}{\|w\|}, \quad d_2 = \frac{|w^T n_n + b|}{\|w\|}$$

$$d_1 = \frac{1}{\|w\|}, \quad d_2 = \frac{1 - 1}{\|w\|} = \frac{1}{\|w\|}$$

$$d = d_1 + d_2 = \frac{2}{\|w\|}$$

Minimize this

Support Vector Machine Objective Function:

Minimize $\frac{1}{2} \|\mathbf{w}\|^2$ where all points have min distance of $1/\|\mathbf{w}\|$.

$\Rightarrow \min \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{Final sum objective}}$ such that $y_i \frac{(\mathbf{w}^T \mathbf{x}_i + b)}{\|\mathbf{w}\|} \geq 1$

$\left. \begin{array}{l} \text{Convex} \\ \text{linear constraint} \end{array} \right\}$

$\nabla_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w}$

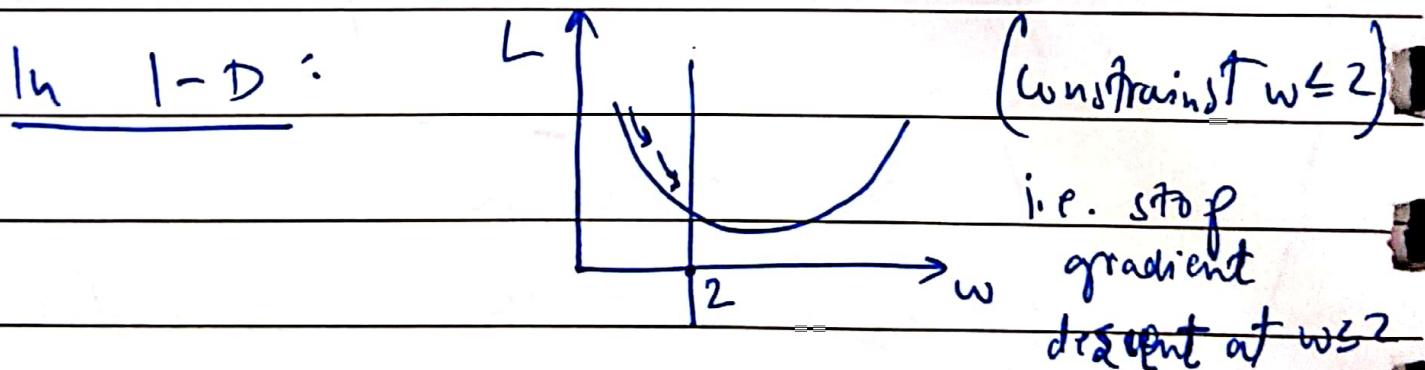
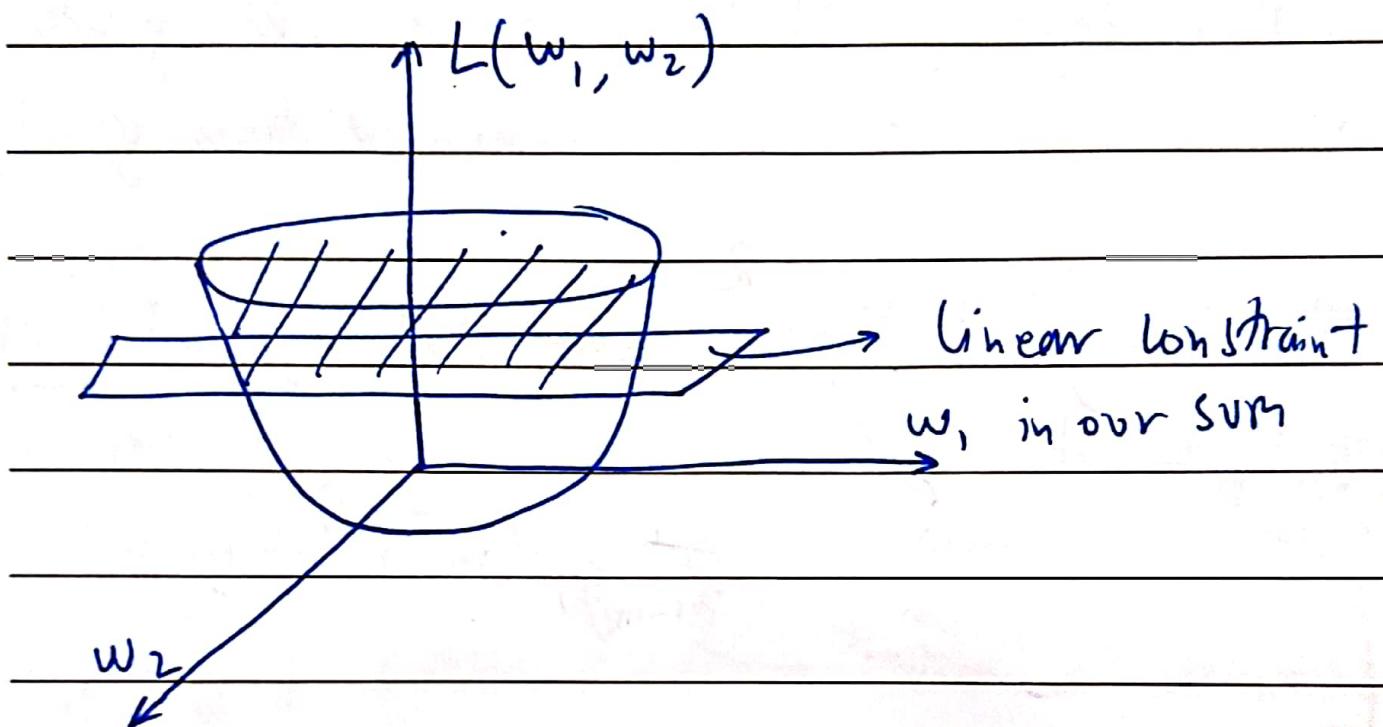
$\nabla_{\mathbf{w}} \mathbf{w}^T \mathbf{w} = \mathbf{w}$ || Derivative

To simplify, we can take square of $\|\mathbf{w}\|$,

This can be solved by :

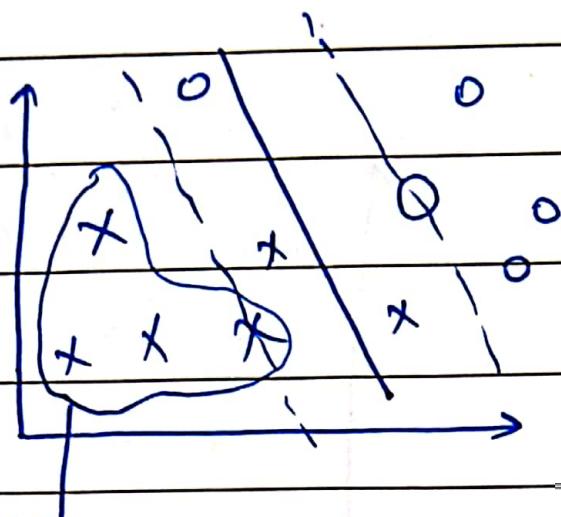
- ① Quadratic Solvers
- ② Lagrangian Duality (Andrew Ng)
- ③ Pegasos (converts unconstrained convex problem into unconstrained problem).

Meaning of convex optimisation problem with some constraint.



Handling outliers in SVM:

- * It is possible for dataset not to be perfectly linearly separable due to presence of outliers.



Allow our objective
to do some errors
some of the training
examples.

$$\text{error } \varepsilon^{(i)} = 0$$

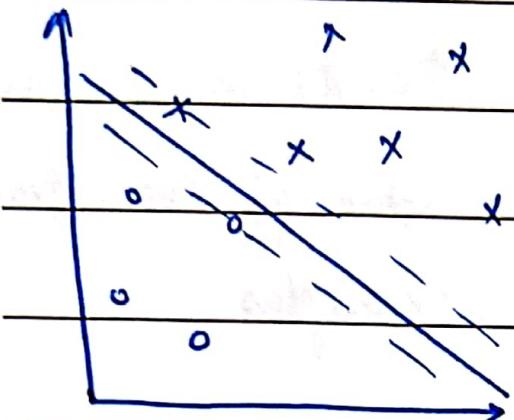
$$(n_i, y_i) \rightarrow \varepsilon^{(i)} \quad // \text{error for each } n^{(i)}$$

Primal
Objective

$$\left. \begin{aligned} & \frac{1}{2} \mathbf{w} \mathbf{w}^T + C \sum_{i=1}^m \varepsilon^{(i)} = \text{loss} \\ & y^{(i)} + (\mathbf{w}^T \mathbf{n}^{(i)} + b) \geq 1 - \varepsilon^{(i)} \end{aligned} \right\}$$

allow some
error.

- * Some points can also have margin ~~$\varepsilon < 1$~~ < 1
↳ even -ve (outlier on opp side).
- * C controls amount of penalty to give when error occurs.



$C \rightarrow \infty$ (very large)

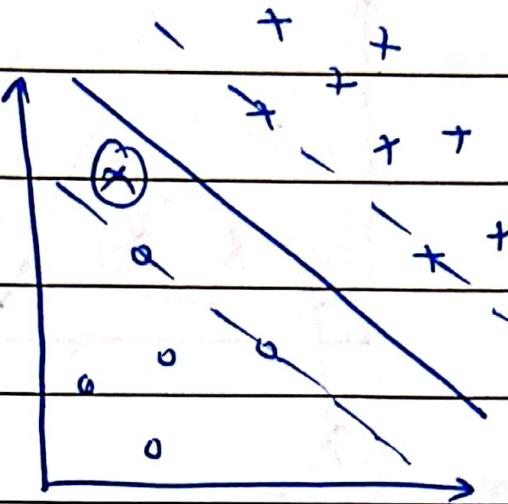
No misclassification as

for loss to be minimum,
all $s^{(i)}$ should be 0 in

this case $\therefore C \leq s^{(i)}$

But margin is less.

* More prone to overfitting.



$C \rightarrow \text{small}$

Some misclassification

for outliers but
margin is max.

Better.

Try different values

of C when using.

Pegasos Algorithm for Unconstrained optimisation.

We did: Non-convex \rightarrow convex with constraints

\downarrow
convex with constraints +
outliers

Now, we convert it is unconstrained convex
optimization so that we can apply gradient
descent, using pegasos.

$$\textcircled{1} \quad \varepsilon^{(i)} \geq 1 - y^{(i)} (\underbrace{w^T x^{(i)} + b}_{t_i}) \quad // \text{from constraint}$$

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



$$\varepsilon_i = 1 - t_i, \text{ if } t_i \leq 1 \\ = 0, \text{ if } t_i > 1$$

If you move away from a (say) positive hyperplane to -ve hyperplane, your error distance will increase, & error will be -ve.

Loss / Sum objective: (Now unconstrained)

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

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

* If point lies in same class & $t_i \uparrow$, then $\varepsilon = 0$.

* $t \leq 1$ when point lies behind the line & here

$t_i \uparrow$ means pt. moving in off direction hence $\varepsilon \uparrow$.

Gradient Descent:

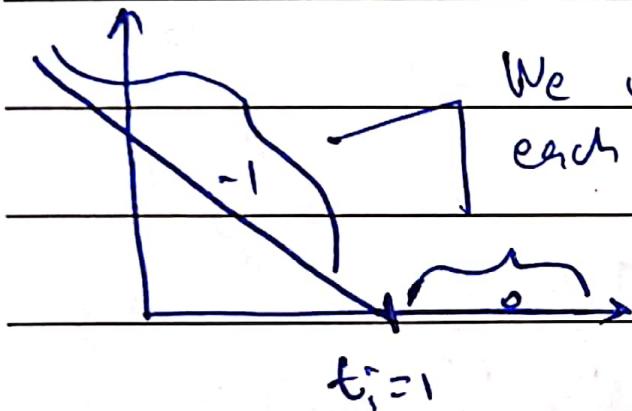
$$w = w - n \nabla_w L$$

~~$b = b - n \nabla_b L$~~

learning rate

$$b = b - n \nabla_b L$$

* func is not differentiable at $t_i=1$.



We will define 2 gradients for each part.

// Also called Hinge-loss

* This technique is called sub-gradient.

$$\nabla_w L = w + c \leq \nabla_w [\max(0, 1-t_i)]$$

$$= w + c \leq \frac{\partial f}{\partial t} \cdot \frac{\partial t}{\partial w} \quad // \text{chain rule}$$

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

$$\nabla_w L = w + c \leq \begin{bmatrix} 0 & \text{if } t_i \geq 1 \\ -1 & \text{if } t_i < 1 \end{bmatrix} \quad \text{y}^{(i)}_m$$

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

$$\nabla_w t_i = y^{(i)} n^{(i)}$$

$$\nabla_b L = 0 + c \begin{bmatrix} 0 & \text{if } t_i \geq 1 \\ -1 & \text{if } t_i < 1 \end{bmatrix} y_i^{(i)}$$

$$\begin{aligned} \nabla_b t_i &= \nabla_b y_i^{(i)} \cdot (w^T y_i^{(i)} + b) \\ &= y_i^{(i)} \end{aligned}$$

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

$$w = w - \eta w + \begin{cases} 0, & t_i \geq 1 \\ \eta c(y_i; n_i), & t_i < 1 \end{cases}$$

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

Final formula for weight update rule.