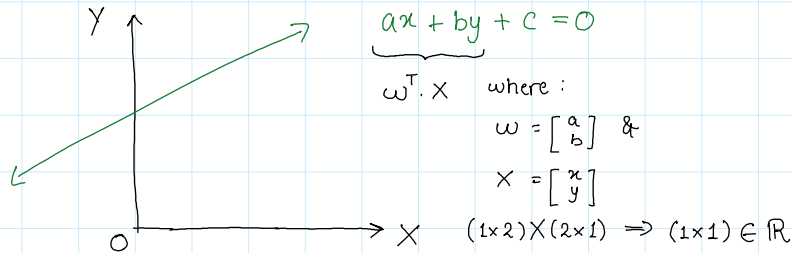


1. In an  $n$ -dimensional space, the equation of a "hyperplane" is given by:

$$\Rightarrow \vec{w}^T \vec{x} + b = 0$$

where  $\vec{w}$  and  $\vec{x}$  are  $n$  dimensional vectors and  $b$  is the bias term.  $\in \mathbb{R}$

Eg: 2-D space (1-D line)



So, we can always characterize or define a  $d$ -dimensional hyperplane with  $\vec{w}$  &  $b$ .  
(Think of  $w$  as the slope of a 1-D line and  $b$  as the intercept.)

Thus, for a line -

$$ax + by + c = 0 \Leftrightarrow w_2 x_1 + w_1 x_0 + w_0 = 0$$

$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix}^T \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} + w_0 = 0$$

for a plane -

$$ax + by + cz + d = 0 \Leftrightarrow w_3 x_2 + w_2 x_1 + w_1 x_0 + w_0 = 0$$

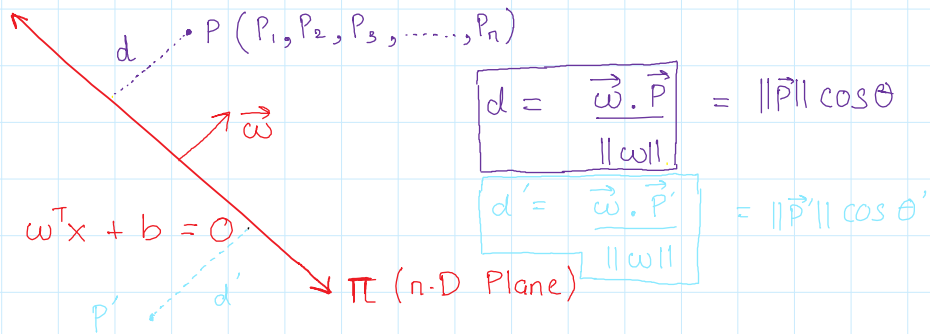
$$\begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}^T \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} + w_0 = 0$$

2. To 'cut' or distribute an  $n$ -dimensional space, we need an  $(n-1)$  dimensional 'hyperplane'.

$\Rightarrow$  Divide a line with a point.

$\Rightarrow$  Cut a cubical space with a plane sheet.

### 3. Distance of a point from a plane (n-D)

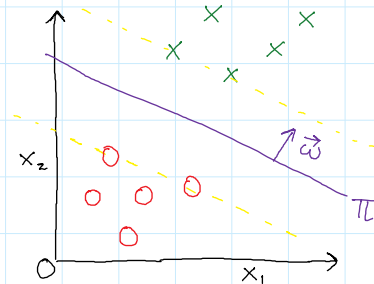


This is pretty intuitive, we take the projection of  $\vec{P}$  on the normal  $\vec{w}$  and divide by its norm to get  $d$ .

ARMED WITH THIS KNOWLEDGE, WE CAN  
ATTACK SUPPORT VECTOR MACHINES ?

**AIM :** To create a Maximum Margin Classifier for some linearly separable data. Or else use Kernel Function(s) to create higher dimensional inner products.

**HYPOTHESIS :** Once your model is ready, it ideally should classify all + samples & - samples correctly. So, it should lie such that all points on one side are + samples and all points on the other side are - samples. with the widest gutter width and "cushioning" on both the sides to maximize generalization.



A sample "crosses" the hyperplane if

$$\begin{aligned} \vec{w} \cdot \vec{x} &\geq c \quad (\text{covered properly in my notes}) \\ \Rightarrow \vec{w} \cdot \vec{x} - c &\geq 0 \\ \Rightarrow \vec{w} \cdot \vec{x} + b &\geq 0 \end{aligned}$$

$$\Rightarrow \vec{w} \cdot \vec{x}_+ + b \geq 0 \quad (\text{for all + samples})$$

$$\vec{\omega} \cdot \vec{x}_- + b \leq 0 \quad (\text{for all } - \text{ samples})$$

$$\Rightarrow \begin{aligned} \vec{\omega} \cdot \vec{x}_+ + b &\geq 0 & \text{for } y = +1 \\ \vec{\omega} \cdot \vec{x}_- + b &\leq 0 & \text{for } y = -1 \end{aligned}$$

$$\Rightarrow y_i (\vec{\omega} \cdot \vec{x}_i + b) \geq 0 \quad \forall i \in [1, 2, \dots, m] \quad (\text{samples})$$

Now let's come back to the training phase. We want to have -

$$\begin{aligned} \vec{\omega} \cdot \vec{x}_i + b &\geq +1 & \text{for } + \text{ samples} \\ \vec{\omega} \cdot \vec{x}_i + b &\leq -1 & \text{for } - \text{ samples} \end{aligned}$$

(covered properly in my notes)

A new explanation for the above eq<sup>ns</sup>:

$$\begin{aligned} \text{Distance of a point from a hyperplane} \\ = \frac{\vec{\omega} \cdot \vec{x}}{\|\vec{\omega}\|} \end{aligned}$$

So for + samples,

$$\frac{\vec{\omega} \cdot \vec{x}_+}{\|\vec{\omega}\|} \geq 1 \quad (\text{some margin threshold})$$

for - samples,

$$\frac{\vec{\omega} \cdot \vec{x}_-}{\|\vec{\omega}\|} \leq -1 \quad (\text{some margin threshold})$$

$$\Rightarrow \begin{aligned} &\text{for } - \text{ samples,} \\ &\text{for } + \text{ samples,} \\ &\frac{\vec{\omega} \cdot \vec{x}_-}{\|\vec{\omega}\|} \leq -1 \quad \frac{\vec{\omega} \cdot \vec{x}_+}{\|\vec{\omega}\|} \geq 1 \end{aligned}$$

Because our normal vector  $\vec{\omega}$  is independent of scaling, changing its magnitude won't do anything because it's only meant for directing our hyperplane. (Discussed properly in my notes)

## OPTIMIZATION FOR HARD MARGIN:

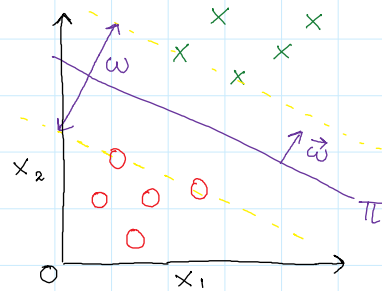
So till now we discussed our expectations from the algorithm in the form of constraints and to sum them up,

$$\omega \cdot x_+ + b \geq +1 \quad \text{①}$$

$$\begin{aligned}\omega \cdot X_+ + b &\geq +1 \\ \omega \cdot X_- + b &\leq -1\end{aligned}$$

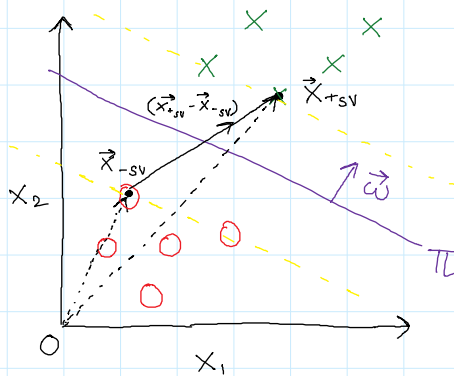
} ①

Our motive was to get the Maximal Margin Classifier which has the widest gutter / street.



Maximize  $\omega$

Note that  $\omega$  is completely different from  $\vec{\omega}$ .



Width of street = Projection of  $(\vec{x}_{+sv} - \vec{x}_{-sv})$  on a unit vector in the direction of  $\vec{\omega}$

$$\begin{aligned}&= (\vec{x}_{+sv} - \vec{x}_{-sv}) \cdot \frac{\vec{\omega}}{\|\vec{\omega}\|} = \frac{\vec{\omega} \cdot \vec{x}_{+sv} - \vec{\omega} \cdot \vec{x}_{-sv}}{\|\vec{\omega}\|} \\ &= \frac{1 - b - (-1 - b)}{\|\vec{\omega}\|} = \boxed{\frac{2}{\|\vec{\omega}\|}} \quad (\text{From ①})\end{aligned}$$

So, we want to maximize  $\frac{2}{\|\vec{\omega}\|}$  or minimize  $\|\vec{\omega}\|$

MINIMIZE  $\left\{ \frac{1}{2} \|\vec{\omega}\|^2 \right\}$  (for mathematical ease)  
(Quadratic problem with just a global minima)

subject to  $y_i \{ \omega \cdot x_i + b \} = 1$  (we don't have inequality

here because our model to solve our optimization problem only depends on SVs)

$$\mathcal{L}(\omega, b) = \frac{1}{2} \omega \cdot \omega - \sum_{i=1}^m \alpha_i [y_i (\omega \cdot x_i + b) - 1]$$

where  $\alpha_i$  :  $i^{\text{th}}$  Lagrangian Multiplier ( $\alpha_i \geq 0$ )

{ More details about Lagrangian are in my NOTES }

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{w}{2} - \sum_{i=1}^m \alpha_i [y_i x_i] = 0$$

$$\Rightarrow \begin{cases} w = \sum_{i=1}^m \alpha_i y_i x_i & \text{--- (2)} \\ \sum_{i=1}^m \alpha_i y_i = 0 & \text{--- (3)} \end{cases}$$

$$\sum_{i=1}^m \alpha_i y_i = 0$$

$\partial b$

Using (2) & (3) in our Primal Lagrangian Problem,

$$\mathcal{L}(w, b) = \frac{1}{2} w \cdot w - \sum_{i=1}^m \alpha_i y_i w \cdot x_i - b \sum_{i=1}^m \alpha_i y_i + \sum_{i=1}^m \alpha_i$$

$$\mathcal{L}(w, b) = \frac{1}{2} w \cdot w - \sum_{i=1}^m \alpha_i y_i w \cdot x_i + \sum_{i=1}^m \alpha_i$$

$$\mathcal{L}(w, b) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) - \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) + \sum_{i=1}^m \alpha_i$$

$$\mathcal{L}_D(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (x_i \cdot x_j)$$

$$\text{where } \alpha_i \geq 0$$

$$\text{and } \sum_{i=1}^m \alpha_i y_i = 0$$

Now, why did we do this?

Well, our primal lagrangian had 3 variables to tune in order to find the minima of the objective (and the lagrangian itself). We converted that to maximizing our new Lagrangian over  $\alpha$ 's by substituting  $w$  and  $b$  as a function of  $\alpha$ . This was done by using the property that derivatives at min = 0.

$$\left( \frac{\partial \mathcal{L}}{\partial w}, \frac{\partial \mathcal{L}}{\partial b} \right)$$

Refer to 'Introduction to Machine Learning: Support Vector Machines' for understanding about Lagrangian.

Eg: minimize  $2 - x^2 - 2y^2$   
 $x, y$

subject to : ①  $x + y - 1 = 0$  (Equality constraint)

② (You can also have inequality constraints)

$$\mathcal{L}(x, y, \alpha) = (2 - x^2 - 2y^2) - \alpha(x + y - 1)$$

and now we have an unconstrained problem with respect to  $x, y$  and  $\alpha$  (Lagrangian mult.)

$$\min \mathcal{L}(x, y, \alpha)$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial x} = 0, \frac{\partial \mathcal{L}}{\partial y} = 0, \frac{\partial \mathcal{L}}{\partial \alpha} = 0$$

$$-2x - \alpha = 0 \quad - \textcircled{1}$$

$$-4y - \alpha = 0 \quad - \textcircled{2}$$

$$x + y - 1 = 0 \quad - \textcircled{3}$$

$$\Rightarrow x + y = 1, \alpha = -2x = -4y \Rightarrow \boxed{x = 2y}$$

$$\boxed{y = 1/3, x = 2/3, \alpha = -4/3}$$

Similarly, you can have  $n$  constraints.

Eg 2:

$$\text{extr.}(f(x, y)) = 8x^2 - 2y$$

$$g(x, y) = x^2 + y^2 - 1 = 0$$

$$\text{extr.}(\mathcal{L}(x, y, \alpha)) = 8x^2 - 2y - \alpha(x^2 + y^2 - 1)$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial x} = 16x - 2\alpha x = 0$$

$$\frac{\partial \mathcal{L}}{\partial y} = -2 - 2y\alpha = 0$$

$$\frac{\partial \mathcal{L}}{\partial \alpha} = -(x^2 + y^2 - 1) = 0$$

$$x^2 + y^2 = 1 \quad - \textcircled{1}$$

$$x(16 - 2\alpha) = 0 \quad - \textcircled{2}$$

$$1 + \alpha y = 0 \quad - \textcircled{3}$$

$$i) x = 0 -$$

$$y^2 = 1 \Rightarrow y = \pm 1$$

$$\alpha = \mp 1$$

$$x = 0, y = 1, \alpha = -1$$

$$x = 0, y = -1, \alpha = +1$$

ii)  $\alpha = 8$ :-

$$y = -1/8 \Rightarrow x = \sqrt{1 - \frac{1}{64}} = \pm \frac{\sqrt{63}}{8}$$

$$x = \pm \frac{\sqrt{63}}{8}, y = -\frac{1}{8}, \alpha = 8$$

Now check which of the 4 solns. give us our maxima/minima.

Once you solve the Dual Problem, you will get your support vectors (those having non-zero  $\alpha_i$ ). After getting your support vectors, you can derive your  $w^*$ .

$$w^* = \sum_{i=1}^n \alpha_i y_i x_i$$

[ This raised a question in my mind, if we used a kernel function to get the  $\alpha$ 's, then we won't be able to use this formula because we won't have the  $x_i$ 's in the higher dimensions. ]

Once we get  $w$ , we can plug in that into:

$$w^* \cdot x_{sv} + b^* = y_{sv} \quad - \textcircled{4}$$

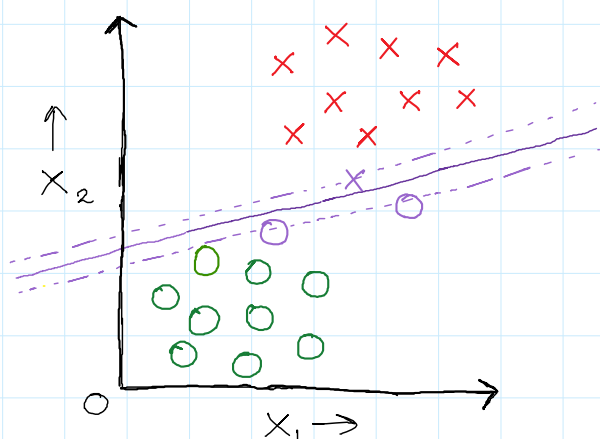
for some SV lying exactly on our margin.

$$\Rightarrow b^* = y_{sv} - w \cdot x_{sv}$$

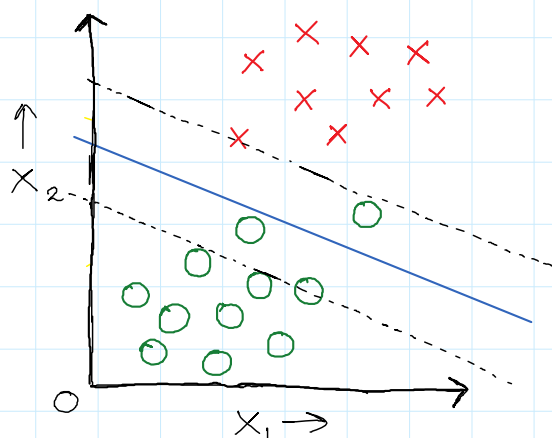
## SOFT MARGIN CLASSIFICATION

For soft margin classifier, we need to allow some misclassifications so that we have a tradeoff between the margin width and the 'correctness' of our classifier. This indicates that we need to have some 'slack' or 'dheel' (दोल) in classifying all the data points correctly so that our margin width does not suffer because of noise in our dataset.

Eg: It is unlikely that any real world data will be linearly separable, even if it is, it might look somewhat like :

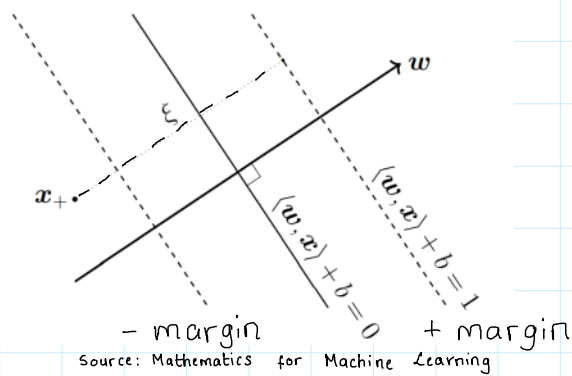


A soft margin SVM will/can allow some slackness in the process which will give us a wider margin at the cost of a few misclassifications.



So we aim at minimizing the total amounts of slacks while maximizing the margin width. Let's introduce a slack variable  $\xi_i$  for every training example  $i = 1, \dots, m$ .





Where  $\epsilon$  quantifies the 'error' or deviation of a + sample  $x_+$  from the +margin:  $w \cdot x + b = +1$ . For a SV,  $\epsilon_{sv} > 1$  means that it has been misclassified as it lies on the other side of the hyperplane. If  $\epsilon_{sv} < 1$ , then it lies inside that margin yet correctly classified.

## OPTIMIZATION:

$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \epsilon_i$$

subject to ①  $y_i(w \cdot x_i + b) \geq (1 - \epsilon_i) \quad \forall i = 1, 2, \dots, m$

The amount of slack for the  $i^{\text{th}}$  sample.

②  $\epsilon_i \geq 0 \quad \forall i = 1, 2, \dots, m$

If  $C$  is large, then even a very little slack will increase the cost. If  $C$  is kept small, then relatively more slack can be allowed.

Thus, reducing  $C$  helps in increasing the no. of SVs inside and on the margin which results in decreasing the variance and reducing over-fitting because our model will generalize better.

New Primal Problem :

New Primal Problem :

$$\mathcal{L}(\omega, b, \xi, \alpha, \mu) = \frac{1}{2} \omega \cdot \omega + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y_i(\omega \cdot x_i + b) - (1 - \xi_i)] - \sum_{i=1}^m \mu_i \xi_i$$

$$\alpha_i \geq 0 \\ \mu_i \geq 0 \quad \forall i = 1, 2, \dots, m$$

$$\min_{\omega, b, \xi} \quad \max_{\alpha, \mu} \quad \mathcal{L}(\omega, b, \xi, \alpha, \mu)$$

NOTE:  $y_i(\omega \cdot x_i + b) \geq 1 - \xi_i$

if  $\xi_i = 0$ , then this  $i^{\text{th}}$  sample is correctly classified.

$$\text{if } y_i(\omega \cdot x_i + b) = 1 - 0$$

then it's a correctly classified S.V.

else if  $\xi_i > 0$ , then this sample might be misclassified.

$$\text{if } y_i(\omega \cdot x_i + b) = 1 - \xi_i \quad \text{and } \xi_i > 1$$

then this 'is' a misclassified data point.

else if  $y_i(\omega \cdot x_i + b) \geq 1 - \xi_i$  is smaller than 1, then it's misclassified.

Partially differentiating the Primal Lagrangian wrt the Primal Variables  $\omega, b$  and  $\xi$  gives us:

$$\frac{\partial \mathcal{L}}{\partial \omega} = \omega - \sum_{i=1}^m \alpha_i y_i x_i$$

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_{i=1}^m \alpha_i y_i$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i} = C - \alpha_i - \mu_i$$

Equating them to zero gives :

$$\boxed{\omega = \sum_{i=1}^m \alpha_i y_i x_i} \quad - \textcircled{1}$$

$$w = \sum_{i=1}^m \alpha_i y_i x_i$$

- ①

$$\sum_{i=1}^m \alpha_i y_i = 0$$

- ②

$$C - \alpha_i - u_i = 0$$

- ③

Substituting ①, ② and ③ after simplifying our Primal Problem :

$$\begin{aligned} \mathcal{L}_P: & \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) - \sum_{i=1}^m \alpha_i y_i \left( \sum_{j=1}^m \alpha_j y_j x_j \cdot x_i \right) - b \sum_{i=1}^m \alpha_i y_i \\ & + \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \alpha_i \xi_i - \sum_{i=1}^m u_i \xi_i + C \sum_{i=1}^m \xi_i \end{aligned}$$

$$\mathcal{L}_D: -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) + \sum_{i=1}^m \alpha_i + \sum_{i=1}^m (C - \alpha_i - u_i) \xi_i$$

$$\mathcal{L}_D(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (x_i \cdot x_j)$$

Clearly,  $\alpha_i \geq 0$  and  $u_i \geq 0 \quad \forall i = 1, 2, \dots, m$  as they are Lagrangian Multipliers.

Given that  $\alpha_i = C - u_i$  (from ③)

$$\Rightarrow 0 \leq \alpha_i \leq C$$

All data points for which  $0 < \alpha_i < C$  holds true, are our Support Vectors that lie exactly on our margin. This is explained by Karush - Kuhn Tucker conditions.

Now, solving the Dual Problem shown above, we will get our Lagrangian Multipliers back. This enables us to get back our Primal Variables  $w$  and  $b$  which parametrize our separating hyperplane.

$$; \dots \dots \dots ;$$

parametrize our separating hyperplane.

$$\omega^* = \sum_{i=1}^m \alpha_i y_i x_i \quad (\text{From } \textcircled{1})$$

Once we get  $\omega$ , we can plug in that into:

$$\omega \cdot x_{sv} + b^* = y_{sv} \quad - \textcircled{4}$$

for some SV lying exactly on our margin.

$$\Rightarrow b^* = y_{sv} - \omega \cdot x_{sv}$$

If everything was hunky-dory, we were done at this point.

BUT

It is possible that none of the SVs lie exactly on the margin i.e. all SVs violate the margin.

Then we can't use  $\textcircled{4}$  as our assumption fails.

By some more math, we can say that in such a case we can find  $b^*$  as follows:

$$b^* = \text{Median}(|y_s - \omega^* \cdot x_s|) \quad \forall s \in SVs.$$