

## Support vector machines.

We have seen that mapping features into another space facilitates the learning of a linear classifier.

Eg: Let  $x = [x_1 \ x_2]$

Define a  $\phi: \mathbb{R}^2 \rightarrow \mathbb{R}^5$ ,

$$z = \phi(x) = [1 \ x_1 \ x_2 \ x_1^2 \ x_2^2 \ x_1 x_2]$$

if  $g(x) = a_0 + a_1 x_1 + a_2 x_2 + a_3 x_1^2 + a_4 x_2^2 + a_5 x_1 x_2$

is a quadratic disc in  $\mathbb{R}^2$ .

$\phi(g(x))$  is a linear disc in  $\phi(x)$  space.

Two major issues with this idea.

1) If we want,  $p^{\text{th}}$  degree poly disc fn in the original feature space ( $\mathbb{R}^m$ ), then the transformed vector,  $z$ , has dim  $O(m^p)$ .

This results in a huge computational cost in both learning & inference.

2) since we need  $O(m^p)$  parameters than  $O(m)$  parameters, ~~we~~ we need large number of examples to achieve generalization.

SVM is designed to offer solution to both of these problems.

Let the training set be

$$\{(x_i, y_i), i=1, \dots, n\}, \quad x_i \in \mathbb{R}^m, \quad y_i \in \{+1, -1\}$$

To begin with assume training set is linearly separable.

$$\exists w \in \mathbb{R}^m \text{ \& } b \in \mathbb{R} \text{ s.t.}$$

$$w^T x_i + b > 0 \quad \forall i \text{ s.t. } y_i = +1$$

$$w^T x_i + b < 0, \quad \forall i \text{ s.t. } y_i = -1$$

$w^T x + b = 0$  is a separating hyperplane

we know that there can be infinitely many separating hyperplanes

$$\Rightarrow \exists \epsilon > 0, \text{ s.t.}$$

$$W^T x_i + b \geq \epsilon, \quad \forall i \text{ s.t. } y_i = +1$$

$$W^T x_i + b \leq -\epsilon \quad \forall i, \text{ s.t. } y_i = -1$$

now,  $w$  can be scaled s.t

$$W^T x_i + b \geq +1 \quad \text{if } y_i = +1$$

$$W^T x_i + b \leq -1 \quad \text{if } y_i = -1$$

equivalently,  $y_i (W^T x_i + b) \geq 1, \forall i$ .

$\Rightarrow$  when training set is separable, any separating hyperplane can be scaled s.t

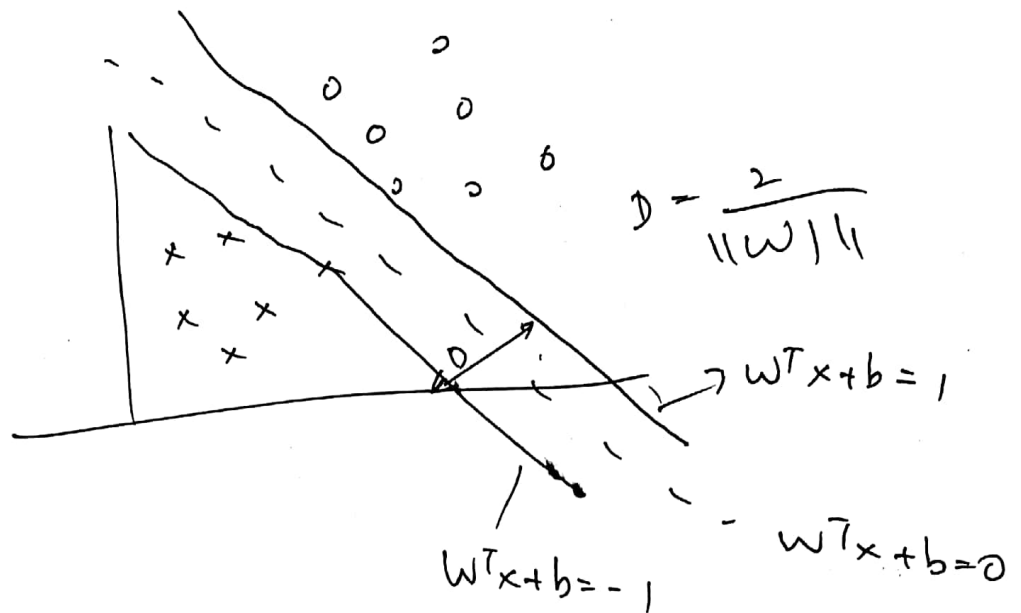
$$y_i (W^T x_i + b) \geq 1, \forall i$$

$\Rightarrow \exists$  no training patterns  $\text{blw}$

$$W^T x + b = +1, \quad \& \quad W^T x + b = -1.$$

Distance b/w these two lines is  $\therefore \frac{2}{\|w\|}$

This is called the margin of the separating hyper plane.



so a good hyperplane is one that maximizes  $D$ .

$\therefore$  Find  $w \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$  to

$$\text{minimize}_{w,b} \quad \frac{1}{2} w^T w$$

$$\text{s.t.} \quad y_i (w^T x_i + b) \geq 1, \quad i=1, \dots, n.$$

## The SVM optimization problem.

Find  $w \in \mathbb{R}^m$ ,  $b \in \mathbb{R}$  s.t

$$\text{minimize } \frac{w^T w}{2}$$

$$\text{subject to } y_i (w^T x_i + b) \geq 1, i=1, \dots, n.$$

Quadratic cost with linear inequality constraints.

## An overview of constrained optimization.

$$\text{minimize } f(x)$$

$$\text{subject to } a_j^T x + b_j \leq 0, j=1, \dots, r.$$

$$f: \mathbb{R}^m \rightarrow \mathbb{R}, a_j \in \mathbb{R}^m, b_j \in \mathbb{R}, j=1, \dots, r.$$

Any point  $x \in \mathbb{R}^m$  is called feasible if

$$a_j^T x + b_j \leq 0, j=1, \dots, r.$$

A constrained optimization problem is solved as follows: Define Lagrangian as follows:

$$L(x, \mu) = f(x) + \sum_{j=1}^r \mu_j (a_j^T x + b_j) \quad (1)$$

If  $f(x)$  is convex,

if  $x_1, x_2 \in \mathbb{R}^n$   $\forall \alpha \in (0, 1)$

$$f(\alpha x_1 + (1-\alpha)x_2) \leq \alpha f(x_1) + (1-\alpha)f(x_2)$$

$$f(x) = x^T x$$

## Kuhn-Tucker Conditions:

An optimization problem with  $f$  convex:

Any  $x^*$  is a global minimum if & only if  $x^*$  is feasible &  $\exists \mu_j^*, j=1, \dots, r$  s.t.

$$1. \nabla_x L(x^*, \mu^*) = 0$$

$$2. \mu_j^* \geq 0, \forall j$$

$$3. \mu_j^* (a_j^T x^* + b_j) = 0, \forall j \rightarrow \text{Complementary slackness}$$

Duality: For a given optimization,

Lagrangian is

$$L(x, \mu) = f(x) + \sum_{j=1}^r \mu_j (a_j^T x + b_j)$$

Here,  $x \in \mathbb{R}^n$  &  $\mu \in \mathbb{R}^r$ .

(171)



Define Dual function :  $q : \mathbb{R}^r \rightarrow [-\infty, \infty]$  by

$$q(\gamma) = \inf_x L(x, \gamma)$$

The dual problem is

$$\text{maximize } q(\gamma)$$

$$\text{subject to } \gamma_j \geq 0, \quad j=1, \dots, r$$

again a constrained opti problem.

over  $\mathbb{R}^r$  &  $\gamma \in \mathbb{R}^r$  are the variables.