

## Primal Dual Relationship.

1. If primal has a soln. so does dual  
& optimal values are equal.

$$q(y^*) = f(x^*)$$

2.  $x^*$  is optimal for primal &  $y^*$  is optimal  
for dual if & only if

a)  $x^*$  is feasible for primal &  $y^*$   
is feasible for dual.

$$b) f(x^*) = L(x^*, y^*) = \min_x L(x, y^*)$$

In the context of SVM,

Primal: minimize  $\frac{1}{2} w^T w$

$$\text{s.t. } 1 - y_i (w^T x_i + b) \leq 0, \quad i = 1, \dots, n.$$

$$L(w, b, \eta) = \frac{1}{2} w^T w + \sum_{i=1}^n \eta_i [1 - y_i (w^T x_i + b)]$$

K.T Cond:

$$\nabla_w L = 0 \Rightarrow w^* = \sum_{i=1}^n \eta_i^* y_i x_i$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{i=1}^n \eta_i^* y_i = 0$$

$$1 - y_i (x_i^T w^* + b^*) \leq 0, \quad \forall i.$$

$$\eta_i^* \geq 0, \quad \eta_i^* [1 - y_i (x_i^T w^* + b^*)] = 0 \quad \forall i$$

Now, if  $S = \{i : \eta_i^* > 0\}$

we have,  $i \in S \Rightarrow y_i (x_i^T w^* + b^*) = 1$

$\Rightarrow x_i$  are the closest to the separating hyperplane.

Thus,  $\{x_i \mid i \in S\}$  are called the "support vectors"

$$\therefore W^* = \sum_{i=1}^n \alpha_i^* y_i x_i = \sum_{i \in S} \alpha_i^* y_i x_i$$

separating hyperplane is a Linear Comb of support vectors

$$W^* = \sum_i \alpha_i^* y_i x_i = \sum_{i \in S} \alpha_i^* y_i x_i$$

$$b^* = y_j - x_j^T W^*, \quad j \text{ s.t. } \alpha_j^* > 0.$$

We want  $M^*$ .

Dual:

$$Q(\gamma) = \inf_{w, b} \left\{ \frac{1}{2} w^T w + \sum_{i=1}^n \gamma_i [1 - \gamma_i (w^T x_i + b)] \right\}$$

We have if  $\sum \gamma_i \gamma_i \neq 0$ , then  $Q(\gamma) = -\infty$

Thus, have another constraint,  $\sum \gamma_i \gamma_i = 0$

Also, min over  $w$  is attained at

$$w = \sum \gamma_i \gamma_i x_i$$

$\therefore w = \sum \gamma_i \gamma_i x_i$  &  $\sum \gamma_i \gamma_i = 0$  in  $Q(\gamma)$

$$Q(\gamma) = \frac{1}{2} w^T w + \sum_{i=1}^n \gamma_i - \sum_{i=1}^n \gamma_i \gamma_i (w^T x_i + b)$$

$$= \sum_i \gamma_i - \frac{1}{2} \sum_i \sum_j \gamma_i \gamma_i \gamma_j \gamma_j x_i^T x_j$$

(176)

∴ The dual is

$$\max_{\mu} \quad q(\mu) = \sum_{i=1}^n \mu_i - \frac{1}{2} \sum_{i,j=1}^n \mu_i \mu_j y_i y_j x_i^T x_j$$

$$\text{s.t.} \quad \mu_i \geq 0, \quad i = 1, \dots, n$$

$$\sum_{i=1}^n y_i \mu_i = 0$$

Observe that in Dual training data only appears as inner products.

Also, optimization is over  $\mathbb{R}^n$ , irrespective of data dimension ( $x_i$ )

$$\Phi(w) = \frac{1}{2} w^T w + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i (w^T x_i + b)$$

Substitute  $w = \sum \alpha_i y_i x_i$ ,  $\sum \alpha_i y_i = 0$ ,

$$\begin{aligned} \Phi(w) &= \frac{1}{2} \left( \sum_i \alpha_i y_i x_i \right)^T \sum_j \alpha_j y_j x_j + \\ &\quad \sum_i \alpha_i - \sum_i \alpha_i y_i x_i^T \left( \sum_j \alpha_j y_j x_j \right) \end{aligned}$$

$$= \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i y_i \alpha_j y_j x_i^T x_j$$

$$\begin{aligned} \text{Dual: } \max_{\alpha} \Phi(\alpha) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \\ \text{s.t. } \alpha_i &\geq 0, \quad i=1, \dots, n, \quad \sum_{i=1}^n y_i \alpha_i = 0. \end{aligned}$$

Final soln for the SVM problem:

$$W^* = \sum \gamma_j^* y_j x_j, \quad b^* = y_j - x_j^T W^*, \quad j: \gamma_j > 0$$

SVM for non-separable case.

i) optimization problem has no feasible point if data is not linearly separable

$$\dots \nexists W, \text{ s.t. } y_i (W^T x_i + b) \geq 1 \quad \forall i$$

However, one can choose a set of slack variables  $\xi_i$  (loose) s.t.

$$y_i (W^T x_i + b) \geq 1 - \xi_i \quad \forall i$$

But the problem now is, every  $W$  is feasible & thus  $\frac{1}{2} W^T W$  will be minimized by  $W=0$ .

To avoid such a degenerative soln, let us modify the cost as follows:

$$\min \frac{1}{2} W^T W + C \sum_{i=1}^n \xi_i$$

$C$  - user defined

$$\text{s.t. } y_i (w^T x_i + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0 \quad \forall i$$

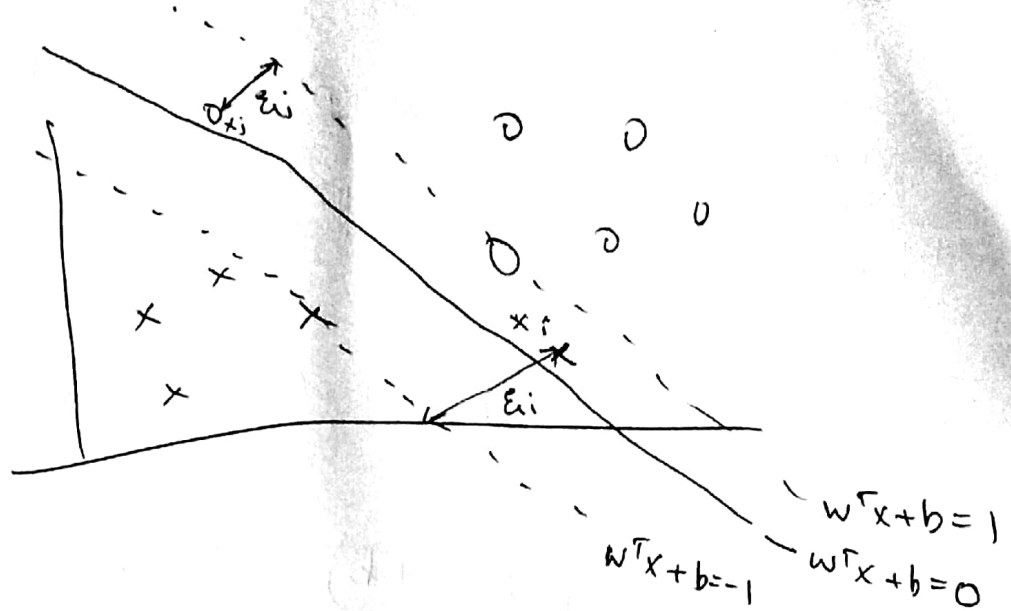
We not only want feasibility but  $\xi_i$  should be as small as possible.

Now, the optimization is over  $w, b$  &  $\xi_i$

Geometrically,  $\xi_i$  measures the extent of violation of optimal separation

i.e., if  $0 < \xi_i < 1$ , then it is a margin error.  
 $\xi_i > 1$ ,  $x_i$  is misclassified.





new opti problem:

$$\min_{W, b, \xi} \frac{1}{2} W^T W + C \sum_{i=1}^n \xi_i$$

$$\text{s.t.} \quad 1 - \xi_i - y_i (W^T x_i + b) \leq 0, \quad i=1, \dots, n$$

$$-\xi_i \leq 0, \quad i=1, \dots, n.$$

Lagrangian:  $L(W, b, \xi, \eta, \alpha) = \frac{1}{2} W^T W + C \sum_{i=1}^n \xi_i$

$$+ \sum_{i=1}^n \eta_i (1 - \xi_i - y_i (W^T x_i + b)) - \sum_{i=1}^n \alpha_i \xi_i$$

$\eta_i \rightarrow$  lagran mul for separability const.

$\alpha_i \rightarrow$  u for  $-\xi_i \leq 0$ .

K. K. T

$$1) \nabla_w L = 0 \Rightarrow w^* = \sum y_i^* y_i x_i$$

$$2) \frac{\partial L}{\partial b} = 0 \Rightarrow \sum y_i^* y_i = 0$$

$$3) \frac{\partial L}{\partial \varepsilon_i} = 0 \Rightarrow y_i^* + d_i^* = c, \forall i$$

$$4) 1 - \varepsilon_i - y_i (w^T x_i + b) \leq 0; \varepsilon_i \geq 0; \forall i$$

$$5) y_i \geq 0; d_i \geq 0, \forall i$$

$$6) y_i (1 - \varepsilon_i - y_i (w^T x_i + b)) = 0; d_i \varepsilon_i = 0, \forall i$$

$w^*$  is same as before.

$$\text{Also, } 0 \leq y_i + d_i = c.$$

Now if  $0 < y_i < c$ , then  $d_i > 0 \Rightarrow \varepsilon_i = 0$ .

$$\Rightarrow \text{from (6), } 1 - y_i (w^T x_i + b) = 0$$

$$\Rightarrow b^* = y_i - x_i^T w^*, \text{ i s.t. } 0 < y_i < c.$$

(183)

To find  $\eta$ , derive ~~Lagra~~ Dual.

$$\Phi(\eta, d) = \inf_{w, b, \xi} L(w, b, \xi, \eta, d)$$

$$L = \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \eta_i (1 - \xi_i - y_i (w^T x_i + b)) - \sum_{i=1}^n d_i \xi_i$$

$L$  has a term :  $\sum_i (C - \eta_i - d_i) \xi_i$

$\Rightarrow$  we need to impose,  $C = \eta_i + d_i \forall i$

If we do, then all terms with  $d_i$  &  $\xi_i$  would vanish & the  $\Phi$  function will be same as before.

only ensure that  $d_i \geq 0$  &  ~~$\eta_i$~~   $\eta_i + d_i = C$ .

in the Dual. which could be easily achieved by  $0 < \eta_i < C$ .

⇒ Dual is

$$\max_{\gamma} \quad \Phi(\gamma) = \sum_{i=1}^n \gamma_i - \frac{1}{2} \sum_{i,j=1}^n \gamma_i \gamma_j \frac{y_i y_j}{x_i^T x_j}$$

$$\text{s.t.} \quad 0 \leq \gamma_i \leq C, \quad i=1..n, \quad \sum_{i=1}^n y_i \gamma_i = 0$$

which is exactly as the previous one  
except for an upper bound on  $\gamma$   
This is another advantage of duality.

∴ SVM in the linear in-separable case:

$$W^* = \sum \gamma_i^* y_i x_i, \quad b^* = y_j - x_j^T W^*,$$

$$\text{s.t.} \quad 0 \leq \gamma_j \leq C.$$

∴ Incorporating slack in the formulation  
ensures that one can find the 'best' hyperplane.