MA 4270    Lecture 8

Last Time : Kernel

— Arise naturally when we consider Regularized cost Func.

$$\theta \in \text{span}\left(\{\phi(x_t)\}\right)_{i}^{n}.$$

Kernel func:  $K: \mathbb{R}^d \times \mathbb{R}^d \longrightarrow \mathbb{R}$    $\exists$ a feature map.  $\phi: \mathbb{R}^d \longrightarrow \mathbb{R}^D$

s.t   $K(x, x') = \langle \phi(x), \phi(x') \rangle$

---

General Convex Optimization & KKT Conditions

Consider        $\min_{x} g_0(x)$    s.t   $g_j(x) \leq 0$   $j=1,\dots, k$    $\rightarrow$ inequality. (k)

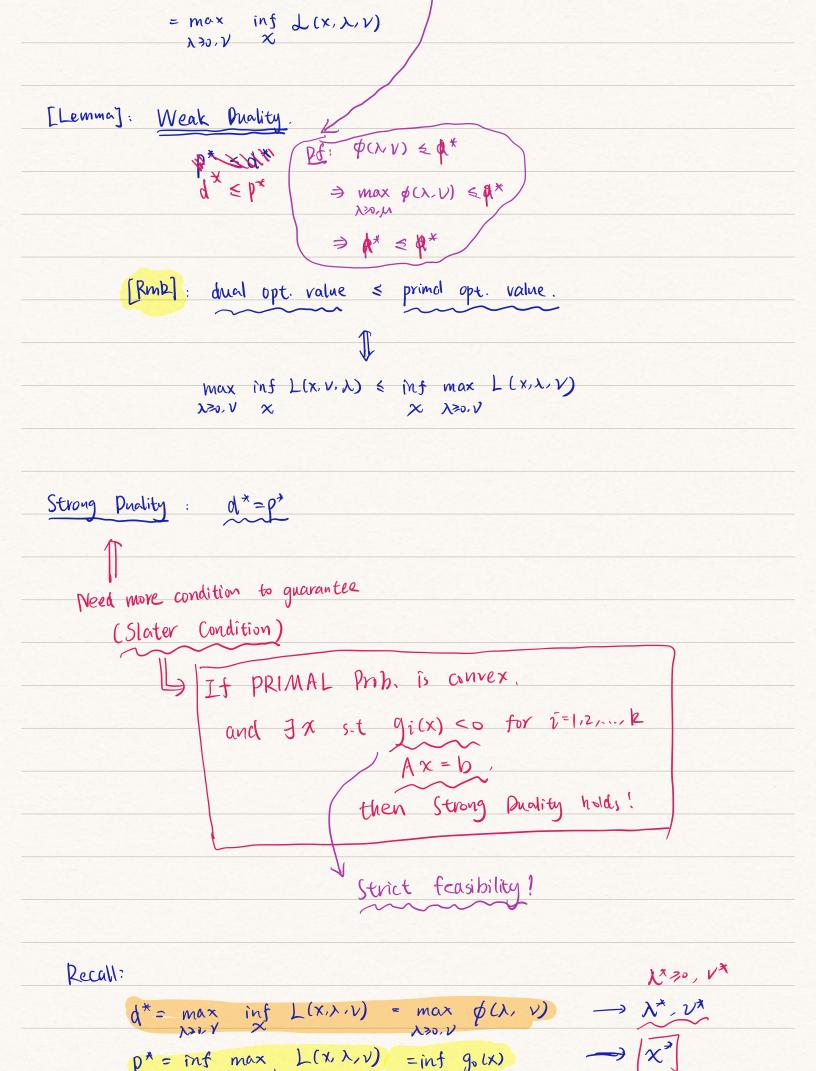$Ax = b$   $\rightarrow$ equality (l)

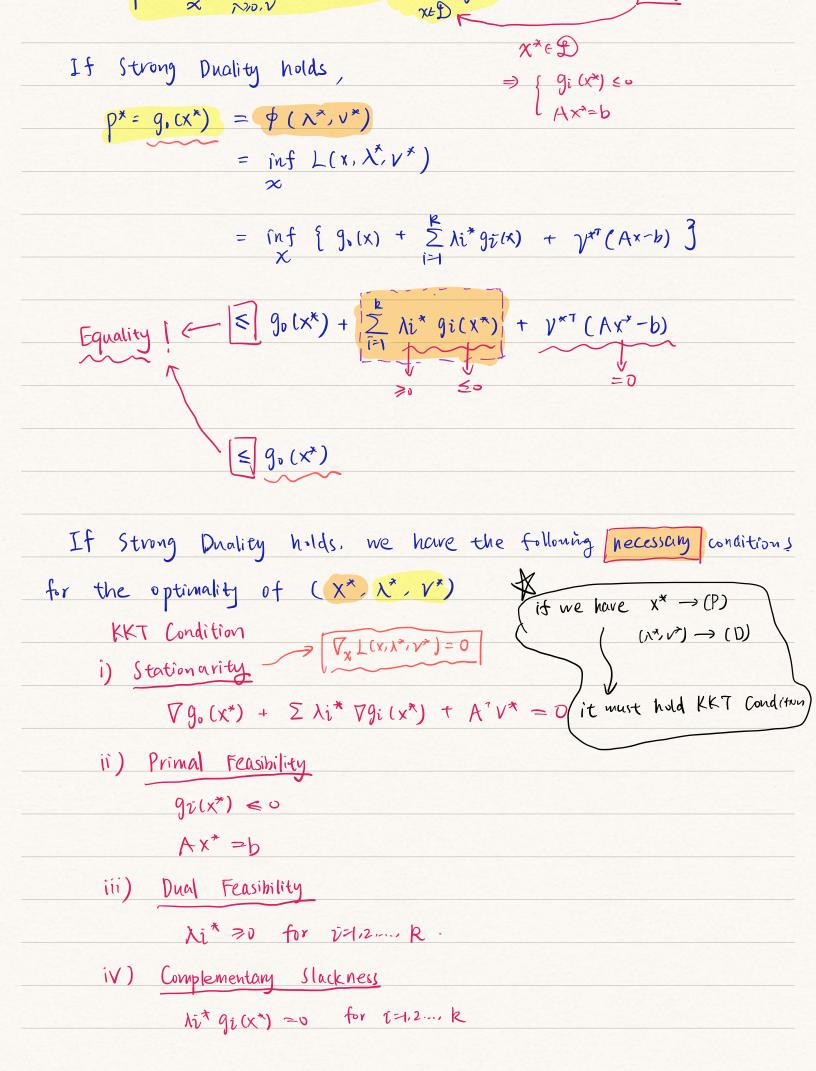Feasible set  $\mathcal{D} = \{ x: g_j(x) \leq 0, j=1,\dots, k . Ax=b\}$

Def: Primal Optimal Value

$$p^* = \inf \{ g_0(x) : x \in \mathcal{D} \}$$

Def .  $\inf \phi = +\infty$   $\Rightarrow$ infeasible.

Def : Lagrangian        $\lambda \in \mathbb{R}^k$    $\nu \in \mathbb{R}^l$.

$$L(x, \lambda, \mu) = g_0(x) + \sum_{i=1}^{k} \lambda_i g_i(x) + \nu^T (Ax - b)$$

| Claim | $p^* = \inf\limits_{x} \left[ \max\limits_{\lambda \geq 0} L(x, \lambda, \nu) \right].$    $\leadsto$ min max Problem |

Pf. ① if $x$ violates some inequl. constraints

then we have $g_{\bar{j}}(x) > 0$ for some $j$

$\quad\hookrightarrow \max\limits_{\lambda \geq 0} L(x, \lambda, \mu) \to +\infty.$ for $x$ violates the constraints

② if $x$ does not violate any constraint

we all have $g_{\bar{j}}(x) \leq 0$   $j=1,\ldots,k$

$\quad\hookrightarrow \max\limits_{\lambda \geq 0} L(x, \lambda, \mu) = g_0(x) \leadsto$ set all $\underline{\lambda}$ to $\underline{0}$

All in all,   $\max\limits_{\lambda \geq 0, \nu} L(x, \lambda, \nu) = \begin{cases} +\infty, & x \notin \mathcal{D} \\ g_0(x), & x \in \mathcal{D} \end{cases}$

Therefore   $p^* = \inf\limits_{x} \max\limits_{\lambda \geq 0, \nu} L(x, \lambda, \nu)$

Def : Lagrangian Dual Function.

$$\phi(\lambda, \nu) = \inf\limits_{x} L(x, \lambda, \nu)$$

Observation (Rmk):   $\boxed{\lambda \geq 0}, \nu, \quad x \in \mathcal{D}$ feasible

$\phi(\lambda, \nu) = \inf\limits_{x} L(x, \lambda, \nu)$

$\quad \leq g_0(x) + \underbrace{\sum\limits_{i=1}^{k} \lambda_i g_i(x)}_{\leq 0} + \underbrace{\nu^T(Ax - b)}_{=0} \quad \forall x$

$\leadsto$ More Condition

$\quad \leq g_0(x) \quad$ for $x \in \mathcal{D}, \lambda_i \geq 0, \nu$

$\Rightarrow \phi(\lambda, \mu) \leq g_0(x) \leq \inf\limits_{x \in \mathcal{D}} g_0(x) = p^*$

Def: Lagrangian Dual Prob.

$$d^* = \max\limits_{\lambda \geq 0, \nu} \phi(\lambda, \nu)$$

$$= \max_{\lambda \geq 0, \nu} \inf_x \mathcal{L}(x, \lambda, \nu)$$

[Lemma]: <u>Weak Duality</u>.

$$p^* \leq d^*$$
$$d^* \leq p^*$$

Pf: $\phi(\lambda, \nu) \leq p^*$

$\Rightarrow \max_{\lambda \geq 0, \mu} \phi(\lambda, \nu) \leq p^*$

$\Rightarrow d^* \leq p^*$

[Rmk]: <u>dual opt. value</u> $\leq$ <u>primal opt. value</u>.

$\Updownarrow$

$$\max_{\lambda \geq 0, \nu} \inf_x L(x, \nu, \lambda) \leq \inf_x \max_{\lambda \geq 0, \nu} L(x, \lambda, \nu)$$

<u>Strong Duality</u> : $\underline{d^* = p^*}$

$\Uparrow$

Need more condition to guarantee

(Slater Condition)

If PRIMAL Prob. is convex,

and $\exists x$ s.t $g_i(x) < 0$ for $i = 1, 2, \ldots, k$

$$Ax = b,$$

then Strong Duality holds!

Strict feasibility!

Recall:

$$d^* = \max_{\lambda \geq 0, \nu} \inf_x L(x, \lambda, \nu) = \max_{\lambda \geq 0, \nu} \phi(\lambda, \nu)$$

$\lambda^* \geq 0, \nu^*$

$\longrightarrow \lambda^*, \nu^*$

$$p^* = \inf_x \max_{\lambda \geq 0, \nu} L(x, \lambda, \nu) = \inf g_0(x)$$

$\longrightarrow \boxed{x^*}$

If Strong Duality holds,

$$p^* = g_0(x^*) = \phi(\lambda^*, \nu^*)$$

$$= \inf_x L(x, \lambda^*, \nu^*)$$

$$= \inf_x \left\{ g_0(x) + \sum_{i=1}^{k} \lambda_i^* g_i(x) + \nu^{*T}(Ax-b) \right\}$$

Equality! $\leftarrow$ $\leq$ $g_0(x^*) + \underbrace{\sum_{i=1}^{k} \lambda_i^* g_i(x^*)}_{} + \underbrace{\nu^{*T}(Ax^*-b)}_{=0}$

$$\geq 0 \qquad \leq 0$$

$$\leq g_0(x^*)$$

$$x^* \in \mathcal{D}$$
$$\Rightarrow \begin{cases} g_i(x^*) \leq 0 \\ Ax^* = b \end{cases}$$

If Strong Duality holds, we have the following **necessary** conditions
for the optimality of $(x^*, \lambda^*, \nu^*)$

KKT Condition

→ $\boxed{\nabla_x L(x, \lambda^*, \nu^*) = 0}$

☆ if we have $x^* \to (P)$
$(\lambda^*, \nu^*) \to (D)$
it must hold KKT Condition

i) Stationarity

$$\nabla g_0(x^*) + \sum \lambda_i^* \nabla g_i(x^*) + A^T \nu^* = 0$$

ii) Primal Feasibility

$$g_i(x^*) \leq 0$$

$$Ax^* = b$$

iii) Dual Feasibility

$$\lambda_i^* \geq 0 \quad \text{for} \quad i=1,2,\dots,k.$$

iv) Complementary Slackness

$$\lambda_i^* g_i(x^*) = 0 \quad \text{for} \quad i=1,2,\dots,k$$

E.g.    SVM without slack

$\min \frac{1}{2}\|\theta\|^2$    s.t    $y_t(\langle\theta, x_t\rangle + \theta_0) \geq 1$    for $t = 1, 2, ..., n$

Strong Duality $\iff$ we can find $(\theta, \theta_0)$ s.t

$$\underline{y_t(\langle x_t, \theta\rangle + \theta_0) \geq 1}$$

$\iff \mathcal{D} = \{(x_t, y_t)\}_{t=1}^n$ is affinely separable!

E.g.    SVM with slack

$\min \frac{1}{2}\|\theta\|^2 + C\sum \xi_t$    s.t    $y_t(\langle x_t, \theta\rangle + \theta_0) \geq 1 - \xi_t$

$$\xi_t \geq 0$$

**Note:**

We can always find $(\theta, \theta_0, \xi)$ s.t

(strict feasibility)

this means Slater Condition always hold

$$y_t(\langle x_t, \theta\rangle + \theta_0) > 1 - \xi_t \quad t = 1, ..., n$$

$$\xi_t > 0 \quad t = 1, 2, ..., n$$

Reason: we can take $\xi_t$ too be BIG enough

Strong Duality always holds!

SVM duality

$\min \frac{1}{2}\|\theta\|^2 + C\sum \xi_t$

s.t    $y_t(\langle x_t, \theta\rangle + \theta_0) \geq 1 - \xi_t$ $\longrightarrow$ (α$_t$)  $\iff$  $1 - \xi_t - y_t(\langle\theta, \phi(x_t)\rangle + \theta_0) \leq 0$

$\xi_t \geq 0$ $\longrightarrow$ (λ$_t$)    $\iff$    $-\xi_t \leq 0$

**Important**

Lagrangian:    $[\underline{\alpha} = (\alpha_1, ..., \alpha_n) \quad \underline{\lambda} = (\lambda_1, ..., \lambda_n)] \longrightarrow$ dual variables

$$\underline{\theta} = (\theta_1, \ldots, \theta_d) \qquad \theta_0 \qquad \underline{\xi} = (\xi_1, \ldots, \xi_n) \rightarrow \text{primal variables}$$

$$\underline{L(\underline{\theta}, \theta_0, \underline{\xi}; \underline{\alpha}, \underline{\lambda})}$$
$$= \tfrac{1}{2}\|\theta\|^2 + C \sum \xi_t + \sum_{t=1}^{n} \boxed{\alpha_t} \underbrace{\left(1 - \xi_t - y_t \left(\langle \underline{\theta}, \Phi(x_t)\rangle + \theta_0\right)\right)}_{g_j} - \sum_{t=1}^{n} \boxed{\lambda_t} \underline{\xi_t}$$

*Lag. v.* (over first box) *Lag. v.* (over last box)

$$\text{primal variable} \qquad \text{dual variable}$$

For $(\underline{\theta^*}, \theta_0^*, \underline{\xi^*}; \underline{\alpha^*}, \underline{\lambda^*})$ to be Primal- Dual Optimal,

we check the KKT ( Necessary ) condition.

① Stationarity.

$$\begin{cases} \dfrac{\partial L}{\partial \underline{\theta}} = 0 \Leftrightarrow \underline{\theta} - \sum \alpha_t y_t \Phi(x_t) = 0 \\[2mm] \dfrac{\partial L}{\partial \theta_0} = 0 \Leftrightarrow \sum \alpha_t y_t = 0 \\[2mm] \dfrac{\partial L}{\partial \xi_t} = 0 \Leftrightarrow C - \alpha_t - \lambda_t = 0 \Leftrightarrow \boxed{\alpha_t + \lambda_t = C} \ \text{For all } t \end{cases}$$

> Actually, all variables here are $(\underline{\theta^*}, \theta_0^*, \underline{\xi^*}; \alpha^*, \lambda^*)$

② Primal Feasibility

$$y_t \left(\langle \underline{\theta}, \Phi(x_t)\rangle + \theta_0\right) \geq 1 - \xi_t \qquad \forall t = 1, 2, \ldots n$$
$$\xi_t \geq 0$$

③ Dual Feasibility

$$\boxed{\alpha_t \geq 0, \qquad \lambda_t \geq 0} \qquad \forall t = 1, 2, \ldots, n .$$

④ Complementary Slackness

$$\alpha_t \left[1 - \xi_t - y_t \left(\langle \underline{\theta}, \Phi(x_t)\rangle + \theta_0\right)\right] = 0 \qquad \forall t = 1, 2, \ldots, n$$
$$\lambda_t \xi_t = 0$$

What is SVs ?

Some Results :

By combining ⬭ , we have : $\alpha_t \in [0, C]$

$\alpha_t$ (diagram: a line from $0$ to $c$)

Partition the $\{(x_t, y_t)\}_{t=1}^{n}$ into 3 disjoint subsets!

① **Non-margin SVs**: $\alpha_t = C > 0$ $\longrightarrow$ CS1

i) $\Rightarrow 1 - \xi_t - y_t (\langle \underline{\theta}, \varphi(x_t) \rangle + \theta_0) = 0$

$\qquad y_t (\langle \underline{\theta}, \varphi(x_t) \rangle + \theta_0) = 1 - \xi_t$

ii) $\Rightarrow \lambda_t = 0 \quad \longleftrightarrow \quad \boxed{\alpha_t + \lambda_t = C}$

$\qquad \Rightarrow \xi_t$ can be positive

② **Margin SVs**: $\alpha_t \in (0, C)$ $\longrightarrow$ CS1.

i) $\Rightarrow 1 - \xi_t - y_t (\langle \underline{\theta}, \varphi(x_t) \rangle + \theta_0) = 0$

$\qquad y_t (\langle \theta, \varphi(x_t) \rangle + \theta_0) = 1 - \xi_t$

ii) $\Rightarrow \lambda_t \neq 0$

$\qquad \Rightarrow \xi_t = 0 \quad \xleftarrow{} \quad$ CS2.

i) + ii) $\Rightarrow y_t (\langle \underline{\theta}, \varphi(x_t) \rangle + \theta_0) = 1$.

$\qquad \qquad \hookrightarrow$ lies in the margin boundary

③ **No SVs**: $\alpha_t = 0$

i) $\lambda_t = C$

$\qquad \Rightarrow \xi_t = 0 \quad \longrightarrow$ CS2.

ii) $\alpha_t = 0 \quad \longrightarrow$ CS1

$\qquad 1 - \xi_t - y_t (\langle \underline{\theta}, \varphi(x_t) \rangle + \theta_0) < 0$

$\qquad y_t (\langle \theta, \varphi(x_t) \rangle + \theta_0) > 1 - \xi_t$

i) + ii) $\qquad y_t (\langle \underline{\theta}, \varphi(x_t) \rangle + \theta_0) > 1$.

$\{x : \langle \varphi(x), \underline{\theta} \rangle + \theta_0 = 0\}$

$$\begin{cases} \bullet \quad \text{Non SVs} \quad \longrightarrow \quad \alpha_t = 0 \\ \bullet \quad \text{Margin SVs} \quad \longrightarrow \quad \alpha_t \in (0, C) \\ \bullet \quad \text{Non-Margin SVs.} \quad \longrightarrow \quad \alpha_t = C \end{cases} \Big\} \text{ SVs.}$$

__Defn__ :   $SV = \{ (X_t, y_t) ; \alpha_t \in (0, C] \}$

__Rmk:__   i) The solution is _sparse_ , i.e., many points have $\alpha_t = 0$

$$\downarrow$$

$$\underline{\text{Non SVs.}}$$

ii) Only the points on margins & those that result in

Margin Errors contribute to the decision of a new

test sample $x'$

$$\Downarrow$$

$\underline{\theta \in \text{span} \{\phi(X_t)\}_{t=1}^{n}} \longrightarrow \boxed{\text{like kernel}}$

$$\hat{y}(x') = \langle \underline{\theta}, \phi(\underline{x'}) \rangle + \theta_0$$

$$= \langle \sum_{t=1}^{n} \alpha_t y_t \phi(X_t), \phi(\underline{x'}) \rangle + \theta_0$$

$$\Downarrow$$

$$\frac{\partial L}{\partial \theta} = 0$$

$$= \sum_{t=1}^{n} \alpha_t y_t K(X_t, \underline{x'}) + \theta_0$$

$\longrightarrow$ High dimension

space

$$= \sum_{t \in SV} \alpha_t y_t K(X_t, X') + \theta_0$$

Offset: $\theta_0$ ? How do we estimate ?

Method: pick a margin SV (i.e. $\alpha_t \in (0, C)$)

$$y_t (< \underline{\theta}, \varphi(X_t) > + \theta_0) = 1$$

$$\Rightarrow y_t (< \sum_{s=1}^{n} \alpha_s y_s \varphi(X_s), \varphi(X_t) > + \theta_0) = 1$$

$$\Rightarrow y_t (\sum_{s=1}^{n} \alpha_s y_s K(X_s, X_t) + \theta_0) = 1$$

$$\Rightarrow \theta_0 = y_t - \sum_{s=1}^{n} \alpha_s y_s K(X_s, X_t)$$

$$= y_t - \sum_{s \in SV} \alpha_s y_s K(X_s, X_t)$$

Question   RBF   $K(x, x') = \exp(-\frac{\beta}{2} \|x - x'\|^2)$ →

| $\beta \uparrow$ | $SV \downarrow$ |
|---|---|
| $\beta \downarrow$ | $SV \uparrow$ |

都不一样

都一样