

MA4270 Lecture 9.

1) Recap of dual SVM formulation

2) Model Selection

3) Structural Risk Management Upper Bound
(SRM)

Primal SVM with slackness & offset.

$$\min_{\underline{\theta}, \theta_0, \underline{\gamma}} \frac{1}{2} \|\underline{\theta}\|^2 + C \sum \gamma_t.$$

$$\text{s.t. } y_t (\langle \underline{\theta}, \underline{\phi}(x_t) \rangle + \theta_0) \geq 1 - \gamma_t$$

$$\gamma_t \geq 0$$

Lagrangian



$$L(\underline{\theta}, \theta_0, \underline{\gamma}; \underline{\alpha}, \underline{\lambda}) = \frac{1}{2} \|\underline{\theta}\|^2 + C \sum \gamma_t + \sum \alpha_t (1 - \gamma_t - y_t (\langle \underline{\theta}, \underline{\phi}(x_t) \rangle + \theta_0)) - \sum \lambda_t \gamma_t$$

KKT Condition :

① Stationarity :

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial \underline{\theta}} = 0 \Rightarrow \underline{\theta} = \sum \alpha_t y_t \underline{\phi}(x_t) \\ \frac{\partial L}{\partial \theta_0} = 0 \Rightarrow \sum \alpha_t y_t = 0 \\ \frac{\partial L}{\partial \gamma_t} = 0 \Rightarrow C = \alpha_t + \lambda_t \end{array} \right. \rightarrow \alpha_t \in [0, C]$$

② feasibility

③ CS

Substitute stationarity conditions into Lagrangian :

$$\phi(\underline{\alpha}, \underline{\lambda}) = \min_{\underline{\alpha}, \underline{\lambda}} L(\underline{\theta}, \theta_0, \underline{\gamma}; \underline{\lambda}, \underline{\alpha})$$

$$= -\frac{1}{2} \sum_{t,s} \alpha_t \alpha_s y_t y_s \langle \underline{\phi}(x_t), \underline{\phi}(x_s) \rangle + \sum_t \alpha_t$$

→ Dual Maximization Prob:

$$\max_{\underline{\alpha}} -\frac{1}{2} \sum_{t,s} \alpha_t \alpha_s y_t y_s \langle \underline{\phi}(x_t), \underline{\phi}(x_s) \rangle + \sum_t \alpha_t$$

$$\text{s.t. : } \alpha_t \in [0, C]$$

$$\sum_t \alpha_t y_t = 0$$

$$\boxed{\frac{\partial L}{\partial \alpha_t} = 0}$$

↓
Stationarity Condition

$$\Leftrightarrow \min_{\underline{\alpha}} \frac{1}{2} \sum_{t,s} \alpha_t \alpha_s y_t y_s \langle \underline{\phi}(x_t), \underline{\phi}(x_s) \rangle - \sum_t \alpha_t$$

$$\text{s.t. } \alpha_t \in [0, C]$$

$$\sum_t \alpha_t y_t = 0$$

$$\Leftrightarrow \boxed{\begin{array}{ll} \min_{\underline{\alpha}} & \frac{1}{2} \underline{\alpha}^T H \underline{\alpha} - \underline{\alpha}^T \underline{1} \\ \text{s.t.} & \alpha_t \in [0, C] \\ & \underline{y}^T \underline{\alpha} = 0 \end{array}}$$

[Rmk]: (i) (P) is a QP in at least D variables.

$$\min \frac{1}{2} \|\underline{\alpha}\|^2 + C \sum_t \beta_t$$

↓
can be very BIG.

$$\text{s.t. } y_t (\langle \underline{\phi}(x_t), \underline{\phi}(x_t) \rangle + \theta_0) \geq 1 - \gamma_t$$

$$\beta_t \geq 0$$

$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$$

(D) is a QP in only n variables

$$\rightarrow \alpha_1, \dots, \alpha_n \in [0, C]$$

→ small dataset.

Dual is cheaper if $D \gg n$ (modern dataset)

ii) Dual only involves the $\langle \phi(x_t), \phi(x_s) \rangle$

$$K(x_t, x_s) = \langle \phi(x_t), \phi(x_s) \rangle.$$

E.g. If we choose RBF Kernel $K(x, x') = \exp(-\gamma \|x - x'\|^2)$

$\Rightarrow \gamma$: "kernelscale" in Matlab.

① if $\gamma \rightarrow 0^+$, $\lim_{\gamma \rightarrow 0^+} K(x, x') = 1$ for all $x, x' \in \mathbb{R}^d$

$$K = \begin{pmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 1 \end{pmatrix} \rightarrow \text{都一样}$$

check the dual prob:

$$\min \frac{1}{2} \sum \alpha_t \alpha_s y_t y_s - \sum \alpha_t$$

$$\text{s.t.: } \alpha_t \in [0, C], \quad \boxed{\sum \alpha_t y_t = 0}$$

$$\Leftrightarrow \min: - \sum \alpha_t$$

$$\text{s.t.: } \alpha_t \in [0, C]$$

Every data point is SV!

$$\sum \alpha_t y_t = 0$$

If: # of '+' = # of '-'

$\Rightarrow \alpha_t = C$ for every t !

guarantee by $\# \text{ of } '+' = \# \text{ of } '-'$

Model Selection : ⇒ New:

Consider 2 candidate kernel for classification/regression



$$K_1(x, x') = 1 + \langle x, x' \rangle \quad \& \quad K_2(x, x') = (1 + \langle x, x' \rangle)^2$$

K_2 is "richer" than K_1 !



For K_1 , the discriminant func (or decision boundary) are functions such that :

$$\begin{aligned} f(x; \theta) &= \theta^T \underline{\phi}^{(1)}(x) \\ &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots \end{aligned}$$

Eg. In 1 dim $f(x; \theta) = \theta_0 + \theta_1 x_1$

For K_2 , the discriminant func (or decision boundary) are functions such that :

$$\begin{aligned} f(x; \theta) &= \theta^T \underline{\phi}^{(2)}(x) \Rightarrow -\text{共 } \binom{d+1}{2} \text{ 个参数} \\ &= \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d + \dots \\ &\quad + \sqrt{2} \theta_{12} x_1 x_2 + \sqrt{2} \theta_{13} x_1 x_3 + \dots + \theta_{dd} x_d^2 + \dots \end{aligned}$$

Informally, whatever discriminant $f \in \mathcal{F}_1$ can produce, K_2 can also produce.



$$\mathcal{F}_1 \subseteq \mathcal{F}_2.$$

More precisely, $\boxed{\mathcal{F}_i = \{ f(\cdot; \theta) : \theta \in \mathbb{R}^{d_i} \}}$

the disri. functions space of K_i

$\forall f_1 \in \mathcal{F}_1$, there $\exists f_2 \in \mathcal{F}_2$ s.t. $f_1 = f_2$!

choose special $\theta_i \in \mathbb{R}^{d_1} \leftarrow \theta_2 \in \mathbb{R}^{d_2}$

Problem: Select a kernel function from a set $\{K_i : i=1,2,\dots,m\}$

where models are associated to f^{\approx} classes F_i .

$\mathcal{F}_1 \subseteq \mathcal{F}_2 \subseteq \dots \subseteq \mathcal{F}_m$.

Nobody tells us use which F_i

$$\mathcal{D}_n = \{(x_t, y_t)\}_{t=1}^n$$

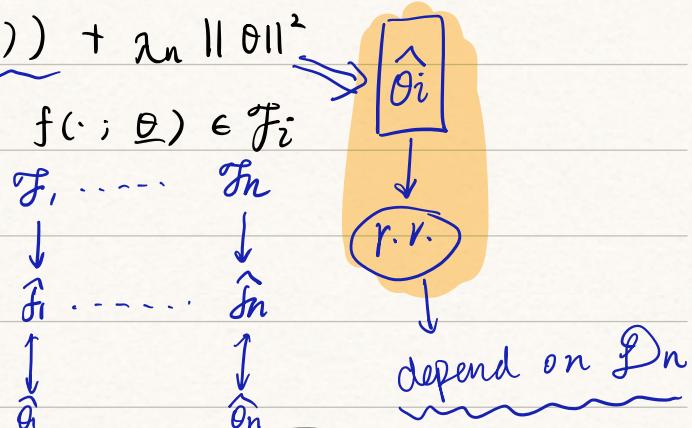
Given that we have to use functions in \mathcal{F}_i , find parameters $\underline{\theta}$ corresponding to $\hat{f}_i \in \mathcal{F}_i$, s.t



Note: this loss should be corresponding to the certain classifier

$$J(\underline{\theta}) = \sum_t \text{loss}(y_t, f(x_t; \underline{\theta})) + \lambda_n \|\underline{\theta}\|^2$$

is minimized over over $\underline{\theta}$ s.t $f(\cdot; \underline{\theta}) \in \mathcal{F}_i$



$$\hat{f}_i(x) = f(x; \hat{\theta}_i) \rightarrow \hat{f}_i \in \mathcal{F}_i.$$

Each parameter of $\underline{\theta}$ has an associated

RISK

Expected Risk

$$\text{RISK} : R(\underline{\theta}) = \mathbb{E}_{(x,y) \sim P} [\text{Loss}_{0-1}(y, f(x; \underline{\theta}))]$$

target our prediction over $f(\cdot; \underline{\theta})$

(x, y) has joint dist. $P(x, y) = P$.

Unknown Prob dist

$$\text{Loss}_{0-1}(y, \bar{y}) = \begin{cases} 1, & y \neq \bar{y} \\ 0, & y = \bar{y} \end{cases}$$

→ if we know $P \Rightarrow$ analytically solve!



But actually we don't know!

→ What we can do?

$$J(\theta) = \sum \text{Loss}(y_t, f(x_t, \theta)) + \lambda_n \|\theta\|^2$$

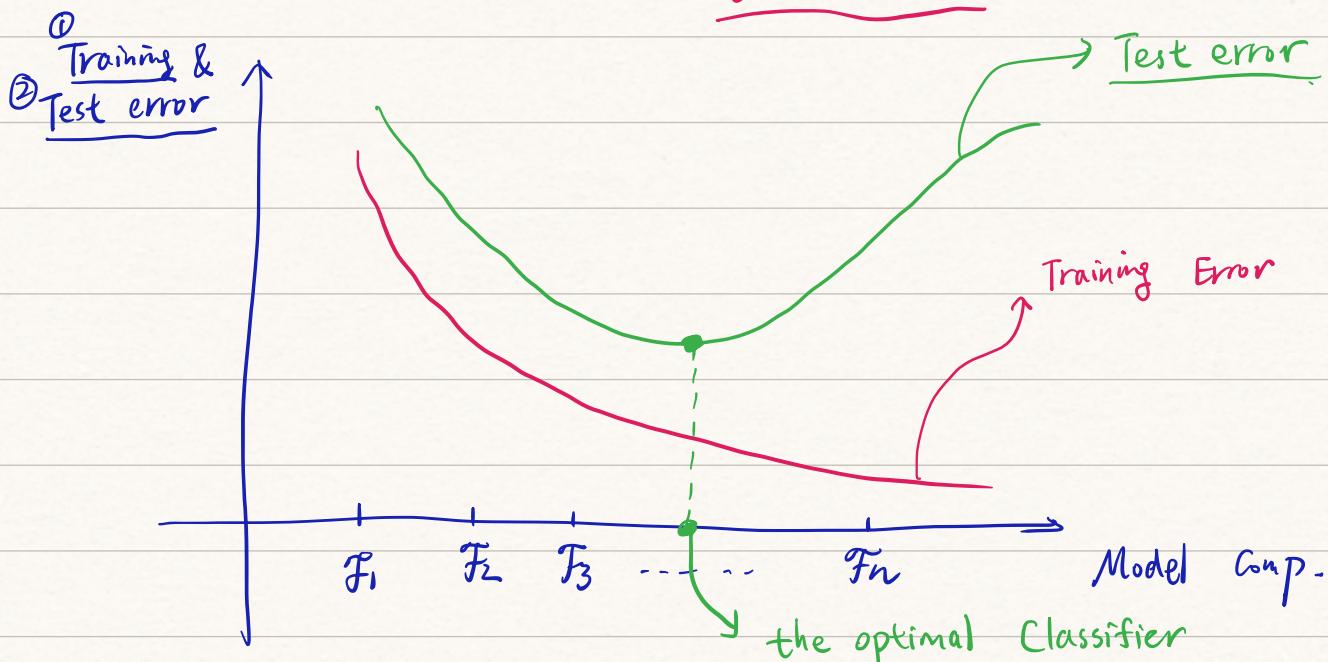
Answer: $\hat{\theta}_i = \underset{\theta \in \mathbb{R}^{d_i}}{\operatorname{argmin}} \{ J(\theta) : f(\cdot; \theta) \in \mathcal{F}_i \}$.

$R(\hat{\theta}) = R(\hat{f}_i)$: still a random variable

risk associated to the function
learned from \bullet

$$(x_t, y_t) \sim P \text{ (i.i.d.)}$$

$\hat{\theta}$ is random



Q: How to select/learn the "sweet spot"?

↓
Correct model order

A: Relate the expected risk $R(\hat{f}_i)$ to the empirical

risk $\hat{R}_n(\hat{f}_i)$

$$R(\hat{f}_i) = \mathbb{E} \left[\underset{y|}{\text{Loss}}(y, \hat{f}_i(x)) \right]$$

$(x,y) \sim P$

expected risk
(unable to calculate)

Unknown

Empirical Risk → obtained from data

$$\hat{R}_n(\hat{f}_i) = \frac{1}{n} \sum_{t=1}^n \underset{y_t}{\text{Loss}}(y_t; \hat{f}_i(x_t))$$

Suppose we can approximate $R(\hat{f}_i)$ → (unknown expected risk) with the empirical risk (known), then we can minimize $\hat{R}_n(\hat{f}_i)$ instead!



Goal: Show that with probability $\geq 1 - \delta$,

$$R(\hat{f}_i) \leq \hat{R}_n(\hat{f}_i) + C(n, \hat{f}_i, \delta)$$

Exp. Risk Emp. Risk.

A uniform Bound

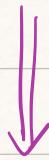
for every $f \in \mathcal{F}$

small!

Use this Upper Bound to select Models

Now perform model selection: → strategy.

$$\hat{i}^* = \operatorname{argmin}_i \{ \hat{R}_n(\hat{f}_i) + \underbrace{C(n, \hat{f}_i, \delta)}_{\downarrow \text{Model Complexity}} \}$$



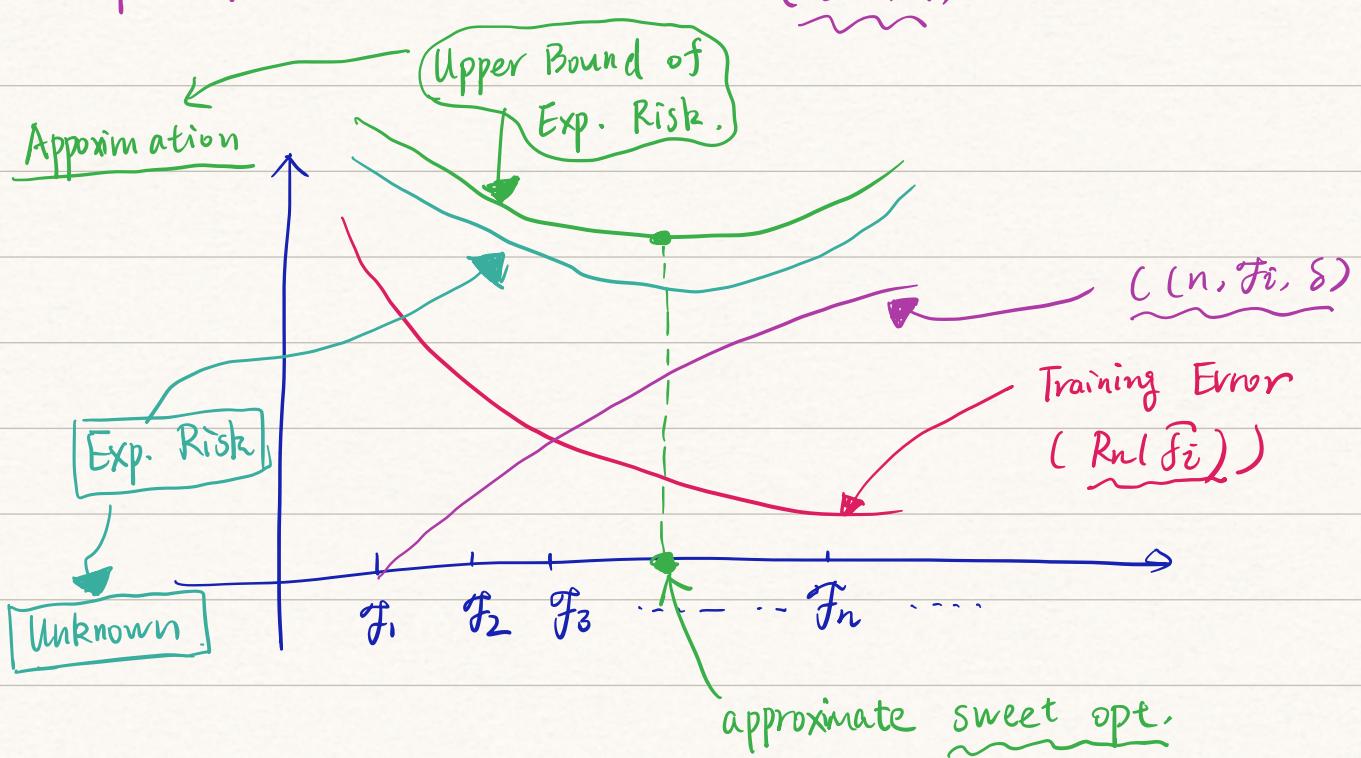
This Principle → Structural Risk Minimization

→ 模型 (SRM)

(SRM)

cf: Minimizing $\hat{R}_n(\hat{f}_i)$ over i is known as

Empirical Risk Minimization (ERM).



Will derive a result for the case in which

$$|\mathcal{F}_i| < \infty$$

$$i=1, 2, \dots, m$$

use VC dimension to extend
this to $|\mathcal{F}_i| = \infty$.

$$\mathcal{F}_i = \{f: \mathbb{R}^d \rightarrow \{\pm 1\} : f(x; \theta) = \text{sgn}(\langle \theta, x \rangle), \theta \in \mathbb{R}^d\}$$

Not a valid func. class because $|\mathcal{F}_i| = \infty$

Looking for a bound s.t.

$$\forall f \in \mathcal{F}_i, \underbrace{R(f)}_{\text{Exp. Risk}} \leq \underbrace{R_n(f)}_{\text{Emp. Risk}} + C(n, \mathcal{F}_i, \delta) \text{ happens}$$

with prob. $\geq 1 - \delta$

Lemma 1: $A_1, \dots, A_k, \Pr(\cup A_i) \leq \sum \Pr(A_i)$

Two Facts!

Lemma 2 : X_1, \dots, X_n i.i.d $E X_1 = \mu$. $\underline{X_i \in [0,1]}$ a.s.

$$\Pr(|\frac{1}{n} \sum X_i - \mu| > \varepsilon) \leq 2 e^{-n\varepsilon^2}$$

Consider : $\Pr(\max_{f \in \mathcal{F}_i} |R(f) - R_n(f)| > \varepsilon) := p$.

this is because $|\mathcal{F}_i| < \infty$

This is the prob. that at least one classifier in \mathcal{F}_i has expected risk deviating from emp. risk by $> \varepsilon$

$$p = \Pr \left(\bigcup_{f \in \mathcal{F}_i} \{ |R(f) - R_n(f)| > \varepsilon \} \right)$$

$$\leq \sum_{f \in \mathcal{F}_i} \Pr(|R(f) - R_n(f)| > \varepsilon) \leftarrow \text{Lemma 1.}$$

$\uparrow \quad \uparrow$
Exp. Risk Emp. Risk.

Associate to each $(x_t, y_t) \sim P$ a r.v.

$$S_t := \prod_{t=1, \dots, n} \{y_t f(x_t) \leq 0\}$$

Note:

① $S_t = 1 \Leftrightarrow f$ makes an error on (x_t, y_t)

② $S_t = 0 \Leftrightarrow f$ makes no error on (x_t, y_t)

$$\Rightarrow R_n(f) = \underbrace{\frac{1}{n} \cdot \sum_t}_{\text{Loss}_0(y_t, f(x_t))} \sum_t S_t$$

$$R(f) = \boxed{\mathbb{E}[S_t]} = \Pr(y_t f(x_t) \leq 0)$$

$$= \mathbb{E}_{(x,y) \sim P} [\text{Loss}_{0-1}(y, f(x))]$$

$$\Pr(|R(f) - R_n(f)| > \varepsilon)$$

$$= \Pr\left(|\mathbb{E}(s_t) - \frac{\sum s_t}{n}| > \varepsilon\right) \quad \underline{s_t \in [0, 1]}$$

$$\leq 2e^{-n\varepsilon^2} \quad \leftarrow \text{Lemma 2}$$

All in all, by Lemma 1 & 2 & $|\mathcal{F}_i| < \infty$

$$\Pr\left(\max_{f \in \mathcal{F}_i} |R(f) - R_n(f)| > \varepsilon\right) \leq |\mathcal{F}_i| 2e^{-n\varepsilon^2} = \delta$$

$$\text{choose } \varepsilon := \sqrt{\frac{\log |\mathcal{F}_i| + \log(\frac{2}{\delta})}{2n}} \Rightarrow \text{with prob } \delta$$

$$\Leftrightarrow \max_{f \in \mathcal{F}_i} |R(f) - R_n(f)| \leq \varepsilon = \sqrt{\frac{\log |\mathcal{F}_i| + \log(\frac{2}{\delta})}{2n}} \quad \text{with prob } \underbrace{1-\delta}_{\sim}$$

Rewrite:

$$|R(f) - R_n(f)| \leq \sqrt{\frac{\log |\mathcal{F}_i| + \log(\frac{2}{\delta})}{2n}} \quad \forall f \in \mathcal{F}_i$$

$$\Rightarrow R(f) \leq R_n(f) + \sqrt{\frac{\log |\mathcal{F}_i| + \log(\frac{2}{\delta})}{2n}} \quad \forall f \in \mathcal{F}_i$$

$$\underbrace{\sqrt{\frac{\log |\mathcal{F}_i| + \log(\frac{2}{\delta})}{2n}}}_{C(n, \mathcal{F}_i, \delta)} \quad \#$$

goes down as $i \rightarrow \infty$ goes up as $i \rightarrow \infty$

Q: $\rightarrow R(f) \rightarrow$ fixed number

$\rightarrow R(\hat{f}_t) \rightarrow$ r.v. ?

$\rightarrow \hat{R}_n(f) \rightarrow$ r.v.

① $R(f) \rightarrow$ exp. risk.

$$E[\text{Loss}_{\text{obj}}(y_t | f(x))]$$

if \hat{f} is attained by:

minimize loss + penalty

Sample dataset

Can we understand as: we do the things in the probability space? \Rightarrow r.v.

\hat{f} ← random r.v.

② $D_n = \{(x_t, y_t)\}_{t=1}^n$ → real world $\rightarrow D = \{(x_t, y_t)\}_{t=1}^n$

theoretical

observation

Emp. risk.

$$\hat{R}_n(f) = \frac{1}{n} \sum \text{Loss}_{\text{obj}}(y_t, f(x_t))$$

can be calculated.

$$f = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{n} \sum \text{loss}(y_t, f(x_t)) + \lambda \|f\|_2^2$$

Conclusion:

Given f , $R(f) \leq R_n(f) + C(n, \mathcal{F}, \delta)$

Compensate for E

$$\begin{aligned}
 R(\hat{f}_i) \Rightarrow & \left\{ \begin{array}{l} \textcircled{1} \quad \mathbb{E} \quad x \\ \textcircled{2} \quad \hat{f}_i \quad x \end{array} \right. \\
 & \rightarrow \hat{R}_n(\hat{f}_i) \Rightarrow \hat{f}_i \quad x \\
 & \downarrow \text{Compensate for } \hat{f}_i \\
 \left\{ \begin{array}{l} \hat{R}_n(\hat{f}_i) \\ \Downarrow \hat{f}_i = \operatorname{argmin} \{ \quad \} \end{array} \right.
 \end{aligned}$$