

① for Non-parametric estimator / parametric estimator

$\hat{\theta} \rightarrow$  estimator of  $\theta$        $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$

$\rightarrow$  MSE decomposition  $\rightarrow$   $\hat{\theta}$  depends on  $X_1, \dots, X_n$

$$\mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)^2]$$

$$= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2 + 0$$

$$= \underbrace{\text{Var}[\hat{\theta}]}_{\text{variance}} + \underbrace{(\mathbb{E}[\hat{\theta}] - \theta)^2}_{\text{bias}^2}$$

$\rightarrow$  similar result holds for Density Estimation

$$\mathbb{E}[(\hat{p}_n(x) - p(x))^2] = \underbrace{(\mathbb{E}[\hat{p}_n(x)] - p(x))^2}_{\text{Bias}^2} + \underbrace{\text{Var}[\hat{p}_n(x)]}_{\text{variance}}$$

② for some Machine Learning Task (Non-parametric Regression included)

Generally speaking, things are more complicated!

$\rightarrow$  Setting

have  $(X, Y)$ , find  $f \in \mathcal{F}$  st  $f(X) \approx Y$

$\rightarrow$  Starting From Decision Theory

a) Assume we have  $(X, Y) \sim P$   $R(f)$

Then consider  $f^* = \underset{f \in \mathcal{F}}{\text{argmin}} \mathbb{E}_{(x, y)} [L(f(x), Y)]$

Assume of Bias & Variance Tradeoff

Specially, if  $L(t_1, t_2) = (t_1 - t_2)^2$   
 then  $\mathbb{E}_{(x,y)} [L(f(x), y)]$

→ Correspond to the CMU Wass...

Lecture 26

Optimal Regression Function

$$= \mathbb{E}_{(x,y)} [(f(x) - y)^2]$$

$$= \mathbb{E}_{(x,y)} [(f(x) - \mathbb{E}[Y|X] + \mathbb{E}[Y|X] - y)^2]$$

$$= \mathbb{E}_{(x,y)} [(f(x) - \mathbb{E}[Y|X])^2] + \mathbb{E}_{(x,y)} [(\mathbb{E}[Y|X] - y)^2]$$

$$+ 2 \mathbb{E}_{(x,y)} [(f(x) - \mathbb{E}[Y|X]) (\mathbb{E}[Y|X] - y)] = 0$$

consider  $\mathbb{E}_y [ \underbrace{(f(x) - \mathbb{E}[Y|X])}_{\text{const}} (\mathbb{E}[Y|X] - y) | X ]$

$$= \text{const} \cdot \underbrace{\mathbb{E}_y [(\mathbb{E}[Y|X] - y) | X]}_0$$

Conclusion

⇒ If we choose  $L(t_1, t_2) = (t_1 - t_2)^2$

$R(f)$

then the optimal model  $f^* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \mathbb{E}_{(x,y)} [L(f(x), y)]$

If we use different  $L(t_1, t_2)$ ,  
 then the corresponding  $f^*(\cdot)$  is different

$$\Leftrightarrow f^*(x) = \mathbb{E}[Y|X=x]$$

↓  
 they are all  
 optimal in different aspects

Also,  $R(f) = \mathbb{E}_{x \sim p_x} [(f(x) - f^*(x))^2] + \mathbb{E}_{(x,y)} [(y - f^*(x))^2]$

MSE

noise (variance)

$$\mathbb{E}_x [\text{Var}[Y|X]]$$

b) After the first step, we consider the randomness

of  $(X, Y) \Rightarrow R(f)$  to achieve the 'optimal model'  $f^*(x)$

↓  
 can be viewed as 'oracle'



We want to use  $\hat{f}_n(x)$  to estimate  $f^*(x)$   
 $\searrow$  RV on  $\mathcal{D}_n$

ERM illustration

[eg.]  $\hat{f}_n = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \hat{R}_n(f)$

where  $\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n L(f(X_i), Y_i)$

$\hat{R}_n(f) \xrightarrow{P} R(f) \quad (X_i, Y_i) \sim P$

Then,  $R(f) = \mathbb{E}_{(x,y)} [L(f(x), y)] \rightarrow$  generalization error

$\downarrow$  useful for Bias & Variance Trade-off  
 $\text{EPE}(\hat{f}_n) = \mathbb{E}_{(x,y,\mathcal{D}_n)} [L(\hat{f}_n(x), y)] \rightarrow$  expected prediction error

- ① Square Loss
- ② use previous result

$\mathbb{E}_{(x,\mathcal{D}_n)} [(\hat{f}_n(x) - f^*(x))^2] + \text{noise}$

$= \mathbb{E}_x \mathbb{E}_{\mathcal{D}_n} [(\hat{f}_n(x) - f^*(x))^2 | X] + \text{noise}$

$= \mathbb{E}_x \left[ \left( \mathbb{E}_{\mathcal{D}_n} [\hat{f}_n(x)] - f^*(x) \right)^2 + \text{Var}_{\mathcal{D}_n} [\hat{f}_n(x)] \right]$

+ noise

$= \mathbb{E}_x [\underbrace{\text{bias}^2(X)}_{\text{w.r.t estimator } \hat{f}_n}] + \mathbb{E}_x [\underbrace{\text{var}(X)}_{\text{w.r.t estimator } \hat{f}_n}] + \underline{\text{noise}}$

$\mathbb{E}_{(x,y)} [(y - f^*(x))^2]$

$\mathbb{E}_x [\text{Var}[Y|X]] \leftarrow$

noise w.r.t  $y$  &  $f^*(x)$

Practically speaking,  $\begin{cases} R(f) \rightarrow \text{infeasible} \\ EPE(\hat{f}_n) \rightarrow \text{infeasible} \end{cases} \rightarrow \text{but can use to explain the } \underline{\text{BIAS \& VARIANCE TRADE-OFF}}$

$\rightarrow$  But can do some kind of simulation:

In our formulation, we require:

$$\underline{y = f(x) + \varepsilon}$$

$$\begin{cases} \mathbb{E}[Y|X=x] = f(x) \leftarrow \mathbb{E}[\varepsilon] = 0 \\ \text{Var}[Y|X=x] = \sigma^2 \text{ for simplicity } \leftarrow \text{Var}[\varepsilon] = \sigma^2 \end{cases}$$

$\rightarrow$  use this to generate data & estimate

$$\begin{cases} \textcircled{1} \text{ bias}(\hat{f}_n(x)) \\ \textcircled{2} \text{ var}(\hat{f}_n(x)) \end{cases}$$

What is feasible?

- 1.  $\hat{R}_n(f) \rightarrow$  unbiased estimator for  $EPE(\hat{f}_n) \Rightarrow$  waste data
- 2.  $CV(f) \rightarrow$  unbiased estimator for  $EPE(\hat{f}_n) \Rightarrow$  computationally expensive

$$EPE(\hat{f}_n) = \mathbb{E}_{(x,y,D)} [L(f(x), y)]$$

Statistical Model

What we are interested?  $\rightarrow$   $\hat{R}_n(f)$

ERM framework

Our interest:  $|R(\hat{f}) - R(f^*)|$

idea is to make sure that  $\hat{f}$  is good enough

$$= | \underbrace{R(\hat{f}) - \hat{R}_n(\hat{f})}_{\textcircled{1}} + \underbrace{\hat{R}_n(\hat{f}) - \hat{R}_n(f^*)}_{\textcircled{2}} + \underbrace{\hat{R}_n(f^*) - R(f^*)}_{\textcircled{3}} |$$

$$= |\textcircled{1}| + |\textcircled{2}| + |\textcircled{3}|$$



→ ② can be controlled since  $\hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \hat{R}_n(f)$

(L function should be good enough)

→ ③ can be controlled due to concentration inequality

$$\mathbb{P}(|\hat{R}_n(f^*) - R(f^*)| \geq \varepsilon) \leq e^{-n \cdot C(\varepsilon)}$$

$$\text{since } \hat{R}_n(f^*) = \frac{1}{n} \sum_{i=1}^n \ell(f^*(x_i), y_i)$$

→ ① Requires uniform bound!

Conceptual idea

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}} |\hat{R}_n(f) - R(f)| \geq \varepsilon\right) \leq C(\varepsilon)$$

$$\|\hat{R}_n - R\|_{f \in \mathcal{F}}$$

{ VC dimension  
Rademacher Theory

More tight bound