

FoSML: Lecture 11

①

• It may be the case that ERM involves NP-hard computation if, for example, we need to do a brute force search over the hypothesis class.

• However, the NP-hardness sometimes results from the non-convexity of the 0-1 loss function. We can address this by developing a convex surrogate loss which upper bounds the 0-1 loss function and has the advantage of being ~~complex~~ convex.

• Let us consider a hypothesis class of \mathbb{R} -valued functions $\mathcal{H} = \{h: \mathcal{X} \rightarrow \mathbb{R}\}$. We can define a binary classifier given h as $\text{sgn}(h)$, i.e. $f_h(x) = \begin{cases} +1, & h(x) \geq 0 \\ -1, & h(x) < 0 \end{cases}$.

• The loss of h at $x \in \mathcal{X}$ is the binary classification error of f_h :

$$\mathbb{1}_{f_h(x) \neq y} = \underbrace{\mathbb{1}_{yh(x) < 0}}_{\substack{y=1 \text{ and } h(x) < 0 \\ \text{or} \\ y=-1 \text{ and } h(x) > 0}} + \underbrace{\mathbb{1}_{\substack{h(x)=0 \\ y=-1}}}_{\substack{\text{boundary} \\ \text{case}}}$$

• The key observation is we can upper bound this as

$$\mathbb{1}_{yh(x) < 0} + \mathbb{1}_{h(x)=0 \wedge y=-1} \leq \mathbb{1}_{yh(x) \leq 0} \quad \star \text{Key inequality!}$$

• We will denote by $R(h)$ the expected error of h :

$$R(h) = \mathbb{E}_{(x,y) \sim D} \left(\mathbb{1}_{f_h(x) \neq y} \right)$$

For any $x \in X$, let $\eta(x)$ denote $\eta(x) = \underbrace{P(y=+1 | x)}_{\substack{\text{stochastic setting, } \\ \neq 1, 0 \text{ in general}}}$. Let D_X denote the marginal distribution over X .

We can decompose the generalization error as $R(h)$

$$\begin{aligned}
 R(h) &= \mathbb{E}_{(x,y) \sim D} (\mathbb{1}_{f_h(x) \neq y}) \\
 &= \mathbb{E}_{\substack{x \sim D_X \\ \cancel{y \sim D}}}} \left[\eta(x) \mathbb{1}_{h(x) < 0} + (1-\eta(x)) \mathbb{1}_{h(x) > 0} + (1-\eta(x)) \hat{\mathbb{1}}_{h(x) = 0} \right] \\
 &= \mathbb{E}_{x \sim D_X} \left[\eta(x) \mathbb{1}_{h(x) < 0} + (1-\eta(x)) \mathbb{1}_{h(x) \geq 0} \right]
 \end{aligned}$$

Minimizing over h , we see the Bayes classifier h^* satisfies

$$h^*(x) = \begin{cases} +1, & \eta(x) \geq \frac{1}{2} \\ -1, & \eta(x) < \frac{1}{2} \end{cases} \Rightarrow h^*(x) = \eta(x) - \frac{1}{2}.$$

Let $R^* = R(h^*)$ denote the Bayes error.

Lemma: For any $h: X \rightarrow \mathbb{R}$ satisfies $R(h) - R^* = 2 \mathbb{E}_{x \sim D_X} \left(|h^*(x)| \mathbb{1}_{h(x)h^*(x) < 0} \right)$.

Proof: For any h , we can write

$$\begin{aligned}
 R(h) &= \mathbb{E}_{x \sim D_X} \left(\eta(x) \mathbb{1}_{h(x) < 0} + (1-\eta(x)) \mathbb{1}_{h(x) \geq 0} \right) \\
 &= \mathbb{E}_{x \sim D_X} \left(\eta(x) \mathbb{1}_{h(x) < 0} + (1-\eta(x)) (1 - \mathbb{1}_{h(x) < 0}) \right)
 \end{aligned}$$

$$= \mathbb{E}_{x \sim D_X} \left(\left[2\eta(x) - 1 \right] \mathbb{1}_{h(x) < 0} + (1 - \eta(x)) \right)$$

$$= \mathbb{E}_{x \sim D_X} \left(2h^+(x) \mathbb{1}_{h(x) < 0} + (1 - \eta(x)) \right)$$

We thus compute $R(h) - R^*$

$$= \mathbb{E}_{x \sim D_X} \left(2h^+(x) \mathbb{1}_{h(x) < 0} + (1 - \eta(x)) \right) - \mathbb{E}_{x \sim D_X} \left(h^+(x) \right)$$

we define $\eta(x) = \frac{1}{2} \mathbb{1}_{h(x) < 0}$

$$= \mathbb{E}_{x \sim D_X} \left(2|h^+(x)| \left(\mathbb{1}_{h(x) \leq 0} - \mathbb{1}_{h^+(x) \leq 0} \right) \right)$$

$$= 2 \mathbb{E}_{x \sim D_X} \left(|h^+(x)| \operatorname{sgn}(h^+(x)) \mathbb{1}_{(h(x)h^+(x) \leq 0) \wedge ((h(x), h^+(x)) \neq (0,0))} \right)$$

$$= 2 \mathbb{E}_{x \sim D_X} \left(|h^+(x)| \mathbb{1}_{h(x)h^+(x) \leq 0} \right)$$

Now, let $\Phi: \mathbb{R} \rightarrow \mathbb{R}$ be convex and non-decreasing. Suppose that for any $u \in \mathbb{R}$, $\mathbb{1}_{u \leq 0} \leq \Phi(-u)$. Define the Φ -loss of $h: \mathcal{X} \rightarrow \mathbb{R}$ at a point $(x, y) \in \mathcal{X} \times \{-1, 1\}$ as $\Phi(-yh(x))$. Then its expected Φ loss is

$$\begin{aligned} \mathcal{L}_\Phi(h) &= \mathbb{E}_{(x,y) \sim D} \left(\Phi(-yh(x)) \right) \\ &= \mathbb{E}_{x \sim D_X} \left(\eta(x) \Phi(-h(x)) + (1 - \eta(x)) \Phi(h(x)) \right). \end{aligned}$$

By construction, $R(h) \leq \mathcal{L}_\Phi(h)$.

For any $x \in \mathcal{X}$, let $u \mapsto L_{\Phi}(x, u)$ be the function defined for $u \in \mathbb{R}$ by $L_{\Phi}(x, u) = \eta(x) \Phi(-u) + (1 - \eta(x)) \Phi(u)$. Then

$$L_{\Phi}(h) = \mathbb{E}_{x \sim D_X} (L_{\Phi}(x, h(x))).$$

Φ is convex $\Rightarrow u \mapsto L_{\Phi}(x, u)$ is convex, being a sum of two convex functions.

Let $h_{\Phi}^{\star} = \mathcal{X} \rightarrow [-\infty, \infty]$ be the Bayes solution for the loss function L_{Φ} .

i.e. $h_{\Phi}^{\star}(x) = \operatorname{argmin}_{u \in [-\infty, \infty]} \eta(x) \Phi(-u) + (1 - \eta(x)) \Phi(u)$.

When $\eta(x) = 0$, $h_{\Phi}^{\star}(x)$ is a minimizer of $\Phi(u)$, and since Φ is non-decreasing, we can set $h_{\Phi}^{\star}(x) = -\infty$. Similarly, when $\eta(x) = 1$, we can set $h_{\Phi}^{\star}(x) = +\infty$.

When $\eta(x) = \frac{1}{2}$, $L_{\Phi}(x, u) = \frac{1}{2} [\Phi(u) + \Phi(-u)]$ so that by convexity,

$L_{\Phi}(x, u) \geq \Phi(0)$, so we can choose $h_{\Phi}^{\star}(x) = 0$ in this case.

In between these extremal values, we may have non-uniqueness of solutions. Pick one arbitrarily.

Let L_{Φ}^{\star} be the Φ -loss of h_{Φ}^{\star} , i.e. $\mathbb{E}_{(x, y) \sim D} (\Phi(-y h_{\Phi}^{\star}(x)))$.

Our goal is the following theorem, which estimates generalization error in terms of L_{Φ}^{\star} :

Theorem: Let Φ be convex and non-decreasing. Suppose $\exists s \geq 1$ and $c > 0$ s.t.
 for all $x \in \mathcal{X}$, $|h^{\Phi}(x)|^s = |y(x) - \frac{1}{2}|^s \leq c^s [L_{\Phi}(x, 0) - L_{\Phi}(x, h^{\Phi}(x))]$.

Then \forall hypotheses h , $R(h) - R^* \leq 2c [L_{\Phi}(h) - L_{\Phi}^*]^{\frac{1}{s}}$.

- Proof next time, but we remark that these assumptions hold broadly:

ex: $\Phi(u) = \max(1+u, 0) \Rightarrow s=1, c=\frac{1}{2}$

$\Phi(u) = \exp(u) \Rightarrow s=2, c=\frac{1}{\sqrt{2}}$

$\Phi(u) = \log_2(1+e^u) \Rightarrow s=2, c=\frac{1}{\sqrt{2}}$
 \vdots