

Fo SML: Lecture 9

• Last time, we showed ERM works well as long as $\hat{R}_S(h) \approx R(h)$. But, in practice we struggle to know $|\hat{R}_S(h) - R(h)|$ in practice.

• An alternative is to consider a family of hypothesis classes $\{H_\gamma\}_{\gamma \in \mathcal{P}}$ with increasing richness. We can then manage the approximation-estimation tradeoff by selecting a "good" choice $\gamma^* \in \mathcal{P}$, not too large and not too small.

• This method is called structural risk minimization, and it works as follows:

Suppose our hypothesis set \mathcal{H} (understood as very large, in most cases) decomposes

as $\mathcal{H} = \bigcup_{k=1}^{\infty} H_k$, where $H_k \subset H_{k+1}$ for all $k \geq 1$. The SRM method

chooses $h_s^{SRM} = \arg \min_{k \geq 1, h \in H_k} \underbrace{R_S(h)}_{\substack{\text{best} \\ \text{empirical} \\ \text{error achievable} \\ \text{in } H_k}} + \underbrace{R_m(H_k)}_{\substack{\text{Rademacher} \\ \text{complexity} \\ \text{of } H_k}} + \underbrace{\sqrt{\frac{\log(k)}{m}}}_{\substack{\text{penalty for working with} \\ \text{a larger hypothesis class} \\ H_k}}$

↳ given a sample of size m

• Let $H_{k(h)}$ be the smallest H_k s.t. $h \in H_k(h)$.

Theorem (SRM learning guarantee): For any $\delta > 0$, with probability exceeding $1 - \delta$, a sample of size m , call it S , is s.t.

$$R(h_s^{SRM}) \leq \inf_{h \in \mathcal{H}} \left(R(h) + 2 \underbrace{R_m(H_{k(h)})}_{\substack{\text{Rademacher} \\ \text{complexity} \\ \text{of } H_{k(h)}}} + \sqrt{\frac{\log(k)}{m}} \right) + \sqrt{\frac{2 \log(\frac{2}{\delta})}{m}}$$

Proof: For notational simplicity, let

$$F_k(h) = \hat{R}_S(h) + \mathcal{R}_m(\mathcal{H}_k) + \sqrt{\frac{\log(k)}{m}}$$

Then by a union bound, we may estimate

$$\begin{aligned} \mathbb{P}\left(\sup_{h \in \mathcal{H}} [R(h) - F_{k(h)}(h)] > \varepsilon\right) &= \mathbb{P}\left(\sup_{k \geq 1} \sup_{h \in \mathcal{H}_k} [R(h) - F_k(h)] > \varepsilon\right) \\ &\leq \sum_{k=1}^{\infty} \mathbb{P}\left(\sup_{h \in \mathcal{H}_k} [R(h) - F_k(h)] > \varepsilon\right) \\ &= \sum_{k=1}^{\infty} \mathbb{P}\left(\sup_{h \in \mathcal{H}_k} [R(h) - \hat{R}_S(h) - \mathcal{R}_m(\mathcal{H}_k)] > \varepsilon + \sqrt{\frac{\log(k)}{m}}\right) \\ &\leq \sum_{k=1}^{\infty} \exp\left(-2m \left[\varepsilon + \sqrt{\frac{\log k}{m}}\right]^2\right) \\ &\leq \sum_{k=1}^{\infty} \exp(-2m\varepsilon^2) \exp(-2\log(k)) \\ &= \sum_{k=1}^{\infty} \exp(-2m\varepsilon^2) \frac{1}{k^2} \\ &= \frac{\pi^2}{6} \cdot \exp(-2m\varepsilon^2) \\ &\leq 2 \exp(-2m\varepsilon^2). \end{aligned}$$

Notice moreover that for any two random variables X_1, X_2 , a union bound argument yields

$$P(X_1 + X_2 > \epsilon) \leq P(X_1 > \frac{\epsilon}{2}) + P(X_2 > \frac{\epsilon}{2}). \quad \text{Now, } \forall h \in \mathcal{H}, \quad (3)$$

$$P(R(h_s^{SRM}) - R(h) - 2\mathcal{R}_m(\mathcal{H}_{k(h)}) - \sqrt{\frac{\log K(h)}{m}} > \epsilon)$$

$$\leq P(R(h_s^{SRM}) - F_{k(h_s^{SRM})}(h_s^{SRM}) > \frac{\epsilon}{2})$$

$$+ P(F_{k(h_s^{SRM})}(h_s^{SRM}) - \cancel{R(h_s^{SRM})} R(h) - 2\mathcal{R}_m(\mathcal{H}_{k(h)}) - \sqrt{\frac{\log K(h)}{m}} > \frac{\epsilon}{2})$$

$$\leq 2\exp(-m\epsilon^2/2) + P(F_{k(h)}(h) - R(h) - 2\mathcal{R}_m(\mathcal{H}_{k(h)}) - \sqrt{\frac{\log K(h)}{m}} > \frac{\epsilon}{2})$$

$\stackrel{=}{=} F_k(h) \forall h \in \mathcal{H}$, by definition of h_s^{SRM}

$$= 2\exp(-m\epsilon^2/2) + P(\hat{R}_S(h) - R(h) - \mathcal{R}_m(\mathcal{H}_{k(h)}) > \frac{\epsilon}{2})$$

$$\leq 2\exp(-m\epsilon^2/2) + \exp(-m\epsilon^2/2)$$

$$= 3\exp(-m\epsilon^2/2).$$

Setting the RHS = δ gets the job done.

Remarks: Suppose $\exists h^* \in \mathcal{H}$ s.t. $R(h^*) = \inf_{h \in \mathcal{H}} R(h)$. Then this result implies

that $\forall h \in \mathcal{H}$, with probability exceeding $1-\delta$,

$$R(h_s^{SRM}) \leq R(h^*) + 2\mathcal{R}_m(\mathcal{H}_{k(h^*)}) + \sqrt{\frac{\log(K(h^*))}{m}} + \sqrt{\frac{2\log(3/\delta)}{m}}.$$

This differs from the estimate we would have if we had known the "best" hypothesis class index $k(h^*)$ a priori... The only additional penalty is the

$$\sqrt{\frac{\log(K(h^*))}{m}} \quad \text{term.}$$

• On the other hand, we may not be able to write $H = \bigcup_{k=1}^{\infty} H_k$ with $H_k \subset H_{k+1}$ (4) and $R_m(H_k)$ converging.

• Even worse from a practical standpoint, actually computing the SRM solution,

namely
$$h_S^{\text{SRM}} = \arg \min_{\substack{k \geq 1 \\ h \in H_k}} \hat{R}_S(h) + R_m(H_k) + \sqrt{\frac{\log(k)}{m}}$$

is NP-hard, and an exhaustive search over k must be performed.

• In practice, one of the most common ways of selecting a predictor is to use cross-validation. In this case, we use some of our training labels to select the hypothesis set H_k .

• As before, let $\{H_k\}_{k \geq 1}$ be a family of nested hypothesis classes.

• Given a sample of size m , S , and $\alpha \in (0, 1)$, we partition S into a training set of size $(1-\alpha)m$ and a validation set of size αm .

• For any $k \in \mathbb{Z}_+$, let $h_{S_1, k}^{\text{ERM}}$ be the ERM solution on S_1 using the hypothesis set H_k . Then

$$h_S^{\text{cv}} = \arg \min_{h \in \{h_{S_1, k}^{\text{ERM}}\}_{k \geq 1}} \hat{R}_{S_2}(h),$$

i.e. we pick the ERM solution which predicts best on S_2 .

Proposition . For any $\alpha > 0$ and any $m \geq 1$,

$$P\left(\sup_{k \geq 1} |R(h_{S_1, k}^{ERM}) - \widehat{R}_{S_2}(h_{S_1, k}^{ERM})| > \epsilon + \sqrt{\frac{\log k}{\alpha m}}\right) \leq 4 \exp(-2\alpha m \epsilon^2).$$

Proof : Next time.