

Recall our result from last time, which we now prove.

Proof of Proposition: By the union bound,

$$\begin{aligned} & \mathbb{P}\left(\sup_{k \geq 1} |R(h_{S_1, k}^{ERM}) - \hat{R}_{S_2}(h_{S_1, k}^{ERM})| > \epsilon + \sqrt{\frac{\log k}{\alpha m}}\right) \\ &= \sum_{k=1}^{\infty} \mathbb{P}\left(|R(h_{S_1, k}^{ERM}) - \hat{R}_{S_2}(h_{S_1, k}^{ERM})| > \epsilon + \sqrt{\frac{\log k}{\alpha m}}\right) \\ &= \sum_{k=1}^{\infty} \mathbb{E}\left(\mathbb{P}\left(|R(h_{S_1, k}^{ERM}) - \hat{R}_{S_2}(h_{S_1, k}^{ERM})| > \epsilon + \sqrt{\frac{\log k}{\alpha m}} \mid S_1\right)\right) \quad (\star) \end{aligned}$$

By conditioning on S_1 , we fix $h_{S_1, k}^{ERM}$. Moreover, S_2 is independent from S_1 .

Then by Hoeffding,

$$\mathbb{P}\left(|R(h_{S_1, k}^{ERM}) - \hat{R}_{S_2}(h_{S_1, k}^{ERM})| > \epsilon + \sqrt{\frac{\log k}{\alpha m}} \mid S_1\right)$$

$$\leq 2 \exp(-2\alpha m (\epsilon + \sqrt{\frac{\log k}{\alpha m}})^2)$$

$$\leq 2 \exp(-2\alpha m \epsilon^2 m - 2 \log k)$$

$$= \frac{2 \exp(-2\alpha m \epsilon^2)}{k^2}$$

Now, summing over k in (\star) yields

$$P\left(\sup_{k \geq 1} |R(h_{S_1, k}^{ERM}) - \hat{R}_{S_2}(h_{S_1, k}^{ERM})| > \varepsilon + \sqrt{\frac{\log k}{2m}}\right)$$

$$\leq \frac{\pi^2}{3} \exp(-2\alpha m \varepsilon^2)$$

$$\leq 4 \exp(-2\alpha m \varepsilon^2), \blacksquare$$

• Our main theoretical result for cross validation is to show that the estimator learned from cross validation is not much worse than what you get from SRM. This is encouraging, because SRM is great in theory, but hard to implement in practice. CV is much easier to implement practically.

Theorem (CV compared to SRM): For any $\delta > 0$, the following holds with probability exceeding $1 - \delta$:

$$R(h_S^{CV}) - R(h_{S_1}^{SRM}) \leq 2 \sqrt{\frac{\log(\max\{k(h_S^{CV}), k(h_{S_1}^{SRM})\})}{\alpha m}} + 2 \sqrt{\frac{\log(4/\delta)}{2\alpha m}}$$

Proof: Note that

$$R(h_S^{CV}) \leq \hat{R}_{S_2}(h_S^{CV}) + \sqrt{\frac{\log(k(h_S^{CV}))}{\alpha m}} + \sqrt{\frac{\log(4/\delta)}{2\alpha m}} \text{ holds with}$$

probability exceeding $1 - \delta$ by our above proposition (just solve for ε in $\delta = 4 \exp(-2\alpha m \varepsilon^2)$). Then since h_S^{CV} ~~is~~ minimizes the empirical error on S_2 , we have

S_2 , we have

$$\begin{aligned}
& \hat{R}_{\alpha 2}(h_{S_1}^{CV}) + \sqrt{\frac{\log(K(h_{S_1}^{CV}))}{2\alpha m}} + 2\sqrt{\frac{\log(4/d)}{2\alpha m}} \\
& \leq \hat{R}_{\alpha 2}(h_{S_1}^{SRM}) + \sqrt{\frac{\log(K(h_{S_1}^{CV}))}{\alpha m}} + 2\sqrt{\frac{\log(4/d)}{2\alpha m}} \\
& \leq R(h_{S_1}^{SRM}) + \sqrt{\frac{\log(K(h_{S_1}^{SRM}))}{\alpha m}} + \sqrt{\frac{\log(K(h_{S_1}^{CV}))}{\alpha m}} + 2\sqrt{\frac{\log(4/d)}{2\alpha m}} \\
& \leq R(h_{S_1}^{SRM}) + 2\sqrt{\frac{\log(\max\{K(h_{S_1}^{CV}), K(h_{S_1}^{SRM})\})}{\alpha m}} + 2\sqrt{\frac{\log(4/d)}{2\alpha m}}
\end{aligned}$$

So, as long as $(1-\alpha)$ is large enough so that SRM on $(1-\alpha)m$ points (i.e. on S_1) gives a good result, while α is large enough so that the $\alpha^{-1/2}$ penalties on the RITs don't hurt too much, we got a good result (bound for CV).

In practice, we may not have m large enough to balance these two conditions, so we use n -fold cross-validation, in which parameters for an algorithm are selected by randomly partitioning the data into training and validation, finding parameter error via cross-validation, then ~~re-splitting~~ ^{re-partitioning} the data.

The more times this split is done (i.e. the more fields), the more costly it is. So it is often the case that 5-fold cross-validation is used.

SRM says "pick a ~~good~~ predictor that fits the data well, while also being simple." This is an example of regularization, a ubiquitous method in statistics and machine learning.

Let $\mathcal{H} = \bigcup_{\gamma > 0} \mathcal{H}_\gamma$ be a nested decomposition of hypothesis classes. For example,

$\mathcal{H}_\gamma = \{x \mapsto w \cdot \Phi(x) \mid \|w\| \leq \gamma\}$ for some function Φ mapping into a high dimensional space (This is closely related to kernel methods, which we shall discuss later)

Here, we are looking at an uncountable sum, since $\gamma \in (0, \infty)$, but we can still

consider SRM:
$$\arg \min_{\gamma > 0, h \in \mathcal{H}_\gamma} \left(\hat{R}_S(h) + R_M(\mathcal{H}_\gamma) + \sqrt{\frac{\log \gamma}{\gamma n}} \right)$$

We can consider more general penalties for complexity than $R_M(\mathcal{H}_\gamma) + \sqrt{\frac{\log \gamma}{\gamma n}}$.

Indeed, let $R: \mathcal{H} \rightarrow \mathbb{R}$ be s.t.
$$\arg \min_{\gamma > 0, h \in \mathcal{H}_\gamma} \hat{R}_S(h) + \underbrace{\text{pen}(\gamma, n)}_{\text{arbitrary penalty}}$$

$$= \arg \min_{h \in \mathcal{H}} \hat{R}_S(h) + \lambda R(h)$$

for some $\lambda > 0$. The second formulation should be familiar to those familiar with signal processing, and we will look at several such regularized problems towards the end of the course.

One does have to take care with how to solve regularized statistical learning problems in practice, since the problems are often non-convex, and thus hard to solve computationally.