

FoSML: Lecture 8

A similar result holds when generalization error of 0 is not achievable, i.e. in the non-realizable case. We will not prove this result:

Theorem (lower bound, non-realizable case): Let \mathcal{H} be a hypothesis set with VC-dimension $d > 1$. Then, $\forall m \geq 1$, for any learning algorithm \mathcal{A} , there exists a distribution over $X \times \{0, 1\}$ such that

$$\mathbb{P}_{S \sim D} \left(R_D(h_S) - \inf_{h \in \mathcal{H}} R_D(h) > \sqrt{\frac{d}{320m}} \right) \geq \frac{1}{64}.$$

Q: How should one construct a good class of hypotheses? If \mathcal{H} is too small, we may be unable to fit the data well. If \mathcal{H} is too big, we may generalize poorly, as shown by our generalization bounds which depend crucially on some of size of \mathcal{H} (e.g. $|\mathcal{H}|$, $VC(\mathcal{H})$, $R(\mathcal{H})$, ...).

We can quantify this tradeoff with the bias/variance or approximation/estimation tradeoff.

More precisely, let \mathcal{H} be a family of functions mapping X to $\{-1, +1\}$. For a hypothesis $h \in \mathcal{H}$, we can quantify its generalization performance by comparing to the Bayes error R^* : $R(h) - R^* \geq 0$.

We can decompose this into two terms, reflecting the approximation (dictated by how well H approximates the Bayes classifier) and the estimation (reflected by how well h does compared to the best it can do) aspects of a statistical learning problem:

$$R(h) - R^* = \underbrace{\left(R(h) - \inf_{h \in H} R(h) \right)}_{\text{estimation}} + \underbrace{\left(\inf_{h \in H} R(h) - R^* \right)}_{\text{approximation}}$$

In general, the approximation error is hard to know, since the underlying distribution D on the data is needed. This is something in the tradition of approximation theory and optimal representations. We will discuss the estimation error instead, since it is closely linked to generalization bounds.

A natural approach is to choose a hypothesis h which minimizes the empirical error on a training sample, then hope (or prove) this will also have low generalization error.

Defn: Let H be a hypothesis class and S a sample. Empirical risk minimization computes the hypothesis with minimal training error: $h_S^{ERM} = \arg \min_{h \in H} \hat{R}_S(h)$.

We can relate the quality of ERM to generalization error as follows:

Proposition: For any sample S , the ERM estimator h_s^{ERM} satisfies

$$P(R(h_s^{ERM}) - \inf_{h \in H} R(h) > \epsilon) \leq P(\sup_{h \in H} |R(h) - \hat{R}_S(h)| > \frac{\epsilon}{2})$$

Proof: For any $\epsilon > 0$, $\exists h_\epsilon$ s.t. $R(h_\epsilon) \leq \inf_{h \in H} R(h) + \epsilon$. Then, using

$R_S(h_s^{ERM}) \leq R_S(h_\epsilon)$, we see

$$\begin{aligned}
R(h_s^{ERM}) - \inf_{h \in H} R(h) &= R(h_s^{ERM}) - R(h_\epsilon) + R(h_\epsilon) - \inf_{h \in H} R(h) \\
&\leq R(h_s^{ERM}) - R(h_\epsilon) + \epsilon \\
&= R(h_s^{ERM}) - \hat{R}_S(h_s^{ERM}) + \hat{R}_S(h_s^{ERM}) - R(h_\epsilon) + \epsilon \\
&\leq R(h_s^{ERM}) - \hat{R}_S(h_s^{ERM}) + \hat{R}_S(h_\epsilon) - R(h_\epsilon) + \epsilon \\
&\leq 2 \cdot \sup_{h \in H} |R(h) - \hat{R}_S(h)| + \epsilon
\end{aligned}$$

Sending $\epsilon \rightarrow 0^+$ gives $R(h_s^{ERM}) - \inf_{h \in H} R(h) \leq 2 \cdot \sup_{h \in H} |R(h) - \hat{R}_S(h)|$.

The result follows immediately. \blacksquare

Remark: We can use ~~constant~~ generalization bounds to control

$$P(\sup_{h \in H} |R(h) - \hat{R}_S(h)| > \frac{\epsilon}{2}),$$

thus acquiring bounds on $P(R(h_s^{ERM}) - \inf_{h \in H} R(h) > \epsilon)$, i.e. on the quality of ERM.