

FoSML: Lecture 5

①

- Recall the notion of Rademacher complexity.
- Our goal is to use the Rademacher complexity to bound the gap between $\hat{R}_S(h)$ (empirical classification error) and $R(h)$ (generalization error). Intuitively, we expect the principle of parsimony/Occam's razor to hold, i.e. lower Rademacher complexity \Rightarrow better bounds on generalization error.

⊛
Thm: Let G be a family of functions mapping Z to $[0,1]$. Then, for any $\delta > 0$, with probability at least $1-\delta$, an iid sample S of size m satisfies:

$$\forall g \in G, \quad (a.) \quad E(g(z)) \leq \frac{1}{m} \sum_{i=1}^m g(z_i) + 2R_m(G) + \sqrt{\frac{\log(2/\delta)}{2m}}$$
$$(b.) \quad \mathbb{P}(g(z)) \leq \frac{1}{m} \sum_{i=1}^m g(z_i) + 2\hat{R}_S(G) + 3\sqrt{\frac{\log(2/\delta)}{2m}}$$

This may be understood as a uniform concentration estimate on the mean, i.e. a quantitative law of large numbers. As $m \rightarrow \infty$, $\sqrt{\frac{\log(2/\delta)}{2m}} \rightarrow \frac{1}{\sqrt{m}}$, so we have the expected convergence rate in m . But, the $R_m(G)/R_S(G)$ terms account for the fact that if G is "complex", there may be some unlucky $g \in G$ that are poorly estimated from the random sample S .

To prove this theorem, we shall make use of another concentration inequality, namely McDiarmid's inequality:

Theorem (McDiarmid's Inequality): Let $\{X_i\}_{i=1}^m$ be a set of independent r.v.'s, and $X_i \in \mathcal{X}_i$ (2)
 Suppose $\exists c_1, \dots, c_m > 0$ s.t. $f: \mathcal{X} \rightarrow \mathbb{R}$ satisfies

$$|f(x_1, \dots, x_i, \dots, x_m) - f(x_1, \dots, x_i', \dots, x_m)| \leq c_i, \quad \forall i=1, \dots, m$$

$$\forall x_1, \dots, x_m, x_i' \in \mathcal{X}_i.$$

Let $f(S) = f(X_1, \dots, X_m)$. Then, $\forall \varepsilon > 0$,

$$(a.) \mathbb{P}(f(S) - \mathbb{E}f(S) \geq \varepsilon) \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^m c_i^2}\right).$$

$$(b.) \mathbb{P}(f(S) - \mathbb{E}f(S) \leq -\varepsilon) \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^m c_i^2}\right).$$

• This says that if the range of f is pointwise bounded (quantified by the c_i 's), then $f(S)$ concentrates (exponentially bounded error estimates) around its mean.

Proof of Thm: Let $\hat{\mathbb{E}}_S(g)$ be the empirical average of g over S , i.e. $\hat{\mathbb{E}}_S(g) = \frac{1}{m} \sum_{i=1}^m g(z_i)$.

We shall apply McDiarmid's inequality to the function

$$\Phi(S) = \sup_{g \in \mathcal{G}} (\mathbb{E}(g) - \hat{\mathbb{E}}_S(g)).$$

Let S, S' be two samples differing in exactly one point, WLOG z_m in S and z_m' in S' .

$$\text{Then: } \Phi(S') - \Phi(S)$$

$$= \sup_{g \in \mathcal{G}} (\mathbb{E}(g) - \hat{\mathbb{E}}_{S'}(g)) - \sup_{g \in \mathcal{G}} (\mathbb{E}(g) - \hat{\mathbb{E}}_S(g))$$

$$\leq \sup_{g \in \mathcal{G}} (\hat{\mathbb{E}}_S(g) - \hat{\mathbb{E}}_{S'}(g))$$

$$\begin{aligned}
&= \sup_{g \in G} \left(\frac{1}{m} \sum_{i=1}^m g(z_i) - \frac{1}{m} \sum_{i=1}^m g(z_i') \right) \\
&= \sup_{g \in G} \left(\frac{g(z_m) - g(z_m')}{m} \right) \\
&\leq \frac{1}{m}, \text{ since } g: \mathbb{Z} \rightarrow [0,1].
\end{aligned}$$

By an identical argument, $\Phi(S) - \Phi(S') \leq \frac{1}{m}$, so $|\Phi(S) - \Phi(S')| \leq \frac{1}{m}$. Then,

by McDiarmid, $\Phi(S) \leq \mathbb{E}_S(\Phi(S)) + \sqrt{\frac{\log(1/\delta)}{2m}}$ with probability at least $1 - \frac{\delta}{2}$.

We now estimate $\mathbb{E}_S(\Phi(S))$ in terms of Rademacher complexity:

$$\begin{aligned}
\mathbb{E}_S(\Phi(S)) &= \mathbb{E}_S \left(\sup_{g \in G} (\mathbb{E}(g) - \hat{\mathbb{E}}_S(g)) \right) \\
&= \mathbb{E}_S \left(\sup_{g \in G} \left(\mathbb{E}_{S'}(\hat{\mathbb{E}}_{S'}(g)) - \hat{\mathbb{E}}_S(g) \right) \right)
\end{aligned}$$

$$\leq \mathbb{E}_{S, S'} \left(\sup_{g \in G} \left(\hat{\mathbb{E}}_{S'}(g) - \hat{\mathbb{E}}_S(g) \right) \right)$$

$$= \mathbb{E}_{S, S'} \left(\sup_{g \in G} \frac{1}{m} \sum_{i=1}^m g(z_i') - g(z_i) \right)$$

$$= \mathbb{E}_{\mathcal{Z}, S, S'} \left(\sup_{g \in G} \frac{1}{m} \sum_{i=1}^m \mathcal{Z}_i (g(z_i') - g(z_i)) \right) \text{ for a Rademacher r.v. } \mathcal{Z} = (\mathcal{Z}_1, \dots, \mathcal{Z}_m)$$

$$\leq \mathbb{E}_{\mathcal{Z}, S'} \left(\sup_{g \in G} \frac{1}{m} \sum_{i=1}^m \mathcal{Z}_i g(z_i') \right) + \mathbb{E}_{\mathcal{Z}, S} \left(\sup_{g \in G} \frac{1}{m} \sum_{i=1}^m -\mathcal{Z}_i g(z_i) \right)$$

$$= 2 \mathbb{E}_{z, S} \left(\sup_{g \in G} \frac{1}{m} \sum_{i=1}^m z_i g(z_i) \right)$$

$$= 2 \mathcal{R}_m(G). \quad \text{This gives (a).}$$

To see (b), we again use McDiarmid's inequality to replace $\mathcal{R}_m(G)$ with $\widehat{\mathcal{R}}_S(G)$:

$$\mathcal{R}_m(G) \leq \widehat{\mathcal{R}}_S(G) + \sqrt{\frac{\log(2/\delta)}{2m}}. \quad \text{The result then follows from union bound, together with}$$

$$\mathbb{P}(S) \leq \mathbb{E}_S(\mathbb{P}(S)) + \sqrt{\frac{\log(2/\delta)}{2m}} \quad \text{and (a).} \quad \blacksquare$$

So far, our results have held for a broad class G . To bring things back to the question of binary classification and generalization error, we need a lemma.

Lemma: Let \mathcal{H} be a family of functions taking values in $\{-1, 1\}$. Let G be the associated family of loss functions with the 0-1 loss:

$$G = \left\{ (x, y) \mapsto \mathbb{1}_{h(x) \neq y} \right\}_{h \in \mathcal{H}}.$$

For any sample $S = \{(x_i, y_i)\}_{i=1}^m \subset \mathcal{X} \times \{-1, 1\}$, let S_x denote the projection onto the first coordinate. Then:

$$\widehat{\mathcal{R}}_S(G) = \frac{1}{2} \widehat{\mathcal{R}}_{S_x}(\mathcal{H}).$$

Proof: By definition, $\widehat{\mathcal{R}}_S(G) = \mathbb{E}_z \left(\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m z_i \mathbb{1}_{h(x_i) \neq y_i} \right)$

$$= \mathbb{E}_z \left(\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m z_i \frac{1 - y_i h(x_i)}{2} \right)$$

$$= \frac{1}{2} \mathbb{E}_z \left(\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m -z_i y_i h(x_i) \right)$$

$$\begin{aligned}
&= \frac{1}{2} \mathbb{E}_D \left(\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m z_i y_i h(x_i) \right) \\
&= \frac{1}{2} \mathbb{E}_D \left(\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m z_i h(x_i) \right) \\
&= \frac{1}{2} \widehat{R}_{S_X}(\mathcal{H}). \quad \blacksquare
\end{aligned}$$

- From $\underline{Thm}(\star)$ and this lemma, it immediately follows that

Thm (Rademacher complexity bound) - binary classification: Let \mathcal{H} be a family of functions taking values in $\{-1, 1\}$. Let D be the distribution over the input space \mathcal{X} . Then $\forall \delta > 0$, with probability exceeding $1 - \delta$ a sample S of size m drawn from D

- satisfies:
- (a) $R(h) \leq \widehat{R}_S(h) + \sqrt{\frac{\log(1/\delta)}{2m}}$
 - (b) $R(h) \leq \widehat{R}_S(h) + \widehat{R}_S(\mathcal{H}) + 3\sqrt{\frac{\log(3/\delta)}{2m}}$