

Lecture 21

Last time

- ▷ Ridge Regression
- ▷ Random design
- ▷ Beyond linear features

Today

- ▷ Radamacher complexity
- ▷ Polynomial discrimination

Radamacher Complexity

We have been studying the relationship between

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n l(\theta, z_i) \text{ and } \min_{\theta \in \Theta} \mathbb{E}_z l(\theta, z).$$

In what follows we will try to understand how to bound the uniform distance between these two.

Goal: Consider $x_i \sim P$ taking values on some set X and let \mathcal{F} be a family of functions of the form $f: X \rightarrow \mathbb{R}$. Our aim is to bound

$$\|P_n - P\|_{\mathcal{F}} := \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E} f(x) \right|.$$

Note that unlike before we will not focus on a specific F .

The following measure of complexity will be instrumental.

Def: The Radamacher Complexity of a set $A \subseteq \mathbb{R}^n$ is given by

$$R(A) := \mathbb{E}_{\varepsilon} \sup_{a \in A} |\langle a, \varepsilon \rangle|,$$

Some people define it without I.I.

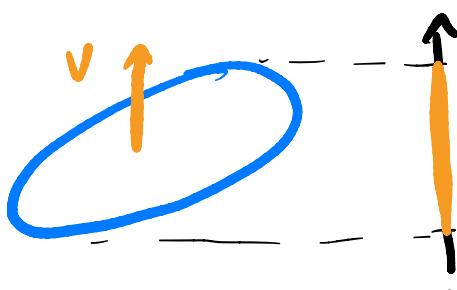
with $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ with ε_i iid Radamacher random variables. \dagger

A geometric interpretation: the support function of a set A is

$$\sigma_A(v) = \sup_{a \in A} \langle v, a \rangle$$

Notice that if v has norm 1, then

"Width of A in the v direction" \downarrow Why? $= \sigma_A(v) + \sigma_{-A}(-v)$



Moreover

$$\mathbb{E} \sigma_A(\varepsilon) \leq R_A \leq 2 \mathbb{E} \sigma_A(\varepsilon).$$

\dagger

Define

$$\mathcal{F}(x) = \{(f(x_1), \dots, f(x_n)) : f \in \mathcal{F}\}.$$

Def: Given a distribution P .
The Radamacher complexity of \mathcal{F} is

$$\begin{aligned} R_n(\mathcal{F}) &:= \mathbb{E}_X R\left(\frac{1}{n} \mathcal{F}(x)\right) \\ &= \mathbb{E}_{X, \varepsilon} \left[\sup \left| \frac{1}{n} \sum \varepsilon_i f(x_i) \right| \right]. \end{aligned}$$

Theorem: For any $n \geq 1$,

$$\mathbb{E} \|P_n - P\|_{\mathcal{F}} \leq 2 R_n(\mathcal{F}).$$

A lower bound also holds (Wainwright Prop 6.11)

Proof: Let $y_1, \dots, y_n \stackrel{\text{iid}}{\sim} P$ and independent of x . Then,

$$\mathbb{E} \|P_n - P\|_{\mathcal{F}}$$

$$= \mathbb{E}_X \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum (f(x_i) - \mathbb{E}_{y_i} f(y_i)) \right| \right]$$

$$\begin{aligned}
&= \mathbb{E}_X \left[\sup_{f \in \mathcal{F}} \left| \mathbb{E}_Y \frac{1}{n} \sum (f(x_i) - f(y_i)) \right| \right] \\
&\stackrel{\text{Jensen's}}{\leq} \mathbb{E}_{X,Y} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum (f(x_i) - f(y_i)) \right| \right] \\
&\stackrel{\text{Symmetric r.v.}}{=} \mathbb{E}_{X,Y,\varepsilon} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum \varepsilon_i (f(x_i) - f(y_i)) \right| \right] \\
&\leq 2 \mathbb{E}_{X,\varepsilon} \left[\sup \left| \frac{1}{n} \sum \varepsilon_i f(x_i) \right| \right] \\
&\leq 2 \mathcal{R}_n(\mathcal{F}). \quad \square
\end{aligned}$$

We also obtain high probability bounds when \mathcal{F} is bounded.

Theorem: Suppose \mathcal{F} satisfies

$$\|f\|_\infty = \sup_{x \in X} |f(x)| \leq b \quad \forall f \in \mathcal{F}.$$

Then, for any $n \geq 1$, $t \geq 0$,

$$\mathbb{P}(\|P_n - P\|_{\mathcal{F}} \leq 2 \mathcal{R}_n(\mathcal{F}) + t) \geq 1 - e^{-\frac{nt^2}{2b^2}}.$$

Proof: This is a consequence of McDiamond's inequality. It suffices to show the bounded difference property with bound $\frac{2b}{n}$.

Set $\bar{f}(x) = f(x) - \mathbb{E}f(x)$. Then

$$\|P_n - P\|_F = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \bar{f}(x_i) \right|.$$

Let x' differ from x only in its i th component. Fix any $\bar{g} \in \mathcal{F}$, then

$$\begin{aligned} & \left| \frac{1}{n} \sum \bar{g}(x_i) \right| - \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum \bar{f}(x'_i) \right| \\ & \leq \left| \frac{1}{n} \sum \bar{g}(x_i) \right| - \left| \frac{1}{n} \sum \bar{g}(x'_i) \right| \\ & \leq \left| \frac{1}{n} (\bar{g}(x_i) - g(x'_i)) \right| \\ & \leq 2b/n, \end{aligned}$$

take sup ^{Gef} and swap the role
of x and x' to establish the
bound. \square

Classes with polynomial discrimination

We transform our problem into
bounding Radamacher complexities.
Next we develop tools to
do so.

Def: \mathcal{F} has polynomial discrimination of order $v \geq 1$ if
 $\forall n \geq 1$ and $\forall x_1, \dots, x_n \in \mathcal{X}$,

$$\#\mathcal{F}(x_1, \dots, x_n) \leq (n+1)^v.$$

Note that even in the simple
case where $\forall f \in \mathcal{F} \quad \text{Im } f = \{1\}$
we can have

$$\#\mathcal{F}(x_1, \dots, x_n) = 2^n$$

Proposition: Suppose \mathcal{F} has polynomial discrimination of order γ . Then

$$R_n(\mathcal{F}) \leq \mathbb{E} \left[\sup_{f \in \mathcal{F}} \sqrt{\frac{1}{n} \sum f(x_i)^2} \right] \sqrt{\frac{2\gamma \log(n+1)}{n}}$$

Proof: The proof relies on the following lemma

Lemma: Suppose $A \subseteq \mathbb{R}^n$

$$R(A) \leq \max_{a \in A} \|a\|_2 \sqrt{2 \log(\#A)}$$

Remember our bound for max of sub-Gaussians. +

Recall that

$$R_n(\mathcal{F}) = \mathbb{E} R\left(\frac{1}{n} \mathcal{F}(X)\right)$$

$$\leq \frac{1}{n} \mathbb{E}_X \max_{f \in \mathcal{F}} \left\| \begin{pmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_n) \end{pmatrix} \right\|_2 \sqrt{2 \log \mathcal{F}(X)}.$$

The final bound follows since \mathcal{F} has polynomial discrimination. \square

In particular if \mathcal{F} is b -bounded then

$$R_n(\mathcal{F}) \leq b \sqrt{\frac{2\gamma \log(n+1)}{n}}.$$

Example: Consider \mathcal{F} to be $\{0, 1\}$ -valued indicator functions of half-intervals

$$f_t(x) = \begin{cases} 1 & \text{if } x \leq t, \\ 0 & \text{otherwise.} \end{cases}$$

Recall the order stats: $x_{(1)} \leq \dots \leq x_{(n)}$,

$$(f_t(x_{(1)}), \dots, f_t(x_{(n)})) = (1, \dots, 1, 0, \dots, 0).$$

Thus, $\# \mathcal{F}(x) \leq n+1$.

Corollary (Glivenko-Cantelli): Let

$$F(t) = P[X \leq t]$$

and \hat{F}_n be the empirical CDF using n iid samples $x_i \sim P$.

Then, for all $\delta \geq 0$,

$$\mathbb{P} \left[\|\hat{F}_n - F\|_{\infty} \geq 4\sqrt{\frac{\log(n+1)}{n}} + \delta \right] \\ \leq e^{-n\delta^2/2}. \quad \blacksquare$$