

## Lecture 23: November 19

Lecturer: Alessandro Rinaldo

Scribe: Ron Yurko

**Note:** *LaTeX template courtesy of UC Berkeley EECS dept.*

**Disclaimer:** *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.*

## 23.1 Last Time

Want to bound supremum of empirical process,

$$\|P_n - P\|_{\mathcal{F}} = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(X_i) - \mathbb{E}[f(X_i)]) \right|$$

where  $(X_1, \dots, X_n) \stackrel{iid}{\sim} P$  on probability space  $(\mathcal{X}, \mathcal{B})$ .  $\mathcal{F}$  is a class of real value functions on  $\mathcal{X}$ , and uniformly bounded:

$$\sup_{x \in \mathcal{X}} |f(x)| \leq B, \quad \forall f \in \mathcal{F}.$$

We will rely on the Rademacher complexity of  $\mathcal{F}$ ,

$$\mathcal{R}_n(\mathcal{F}) = \mathbb{E}_{\tilde{X}, \tilde{\epsilon}} \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \left| \sum_{i=1}^n \epsilon_i f(X_i) \right| \right]$$

**Theorem 23.1** *Let  $\mathcal{F}$  be a class of functions  $(\mathcal{X}, \mathcal{B})$  uniformly bounded by  $B > 0$ . Then for any data generating distribution  $P$  for  $(X_1, \dots, X_n)$  and for all  $t > 0$ ,*

$$\mathbb{P}(\|P_n - P\|_{\mathcal{F}} \leq 2\mathcal{R}_n(\mathcal{F}) + t) \leq \exp\left(\frac{-nt^2}{2B^2}\right)$$

Remark If  $\mathcal{R}_n(\mathcal{F}) \rightarrow 0$  as  $n \rightarrow \infty$ , then  $\|P_n - P\|_{\mathcal{F}} \xrightarrow{a.s.} 0$  (this is actually a if and only if statement).

*Proof.* Proof has 2 parts:

1. Show  $\|P_n - P\|_{\mathcal{F}}$  concentrates around  $\mathbb{E}[\|P_n - P\|_{\mathcal{F}}]$  (easy part)
2. Bound  $\mathbb{E}[\|P_n - P\|_{\mathcal{F}}]$  by  $2\mathcal{R}_n(\mathcal{F})$  (hard part)

Part 1: Use the Bounded Difference inequality  $\rightarrow$  off-the-shelf inequality meant for this situation.

Let  $(X_1, \dots, X_n) \subset \mathcal{X}$  be arbitrary n-tuple of points. Define the function  $G : \mathcal{X}^n \Rightarrow \mathbb{R}$ ,

$$x_1^n = (x_1, \dots, x_n) \Rightarrow \sup_{f \in \mathcal{F}} \frac{1}{n} \left| \sum_{i=1}^n \bar{f}(x_i) \right|,$$

where  $\bar{f}(x_i) = f(x_i) - \mathbb{E}[f(x_i)]$ . Let  $x_1^n, y_1^n$  be in  $\mathcal{X}^n$  such that  $x_i = y_i$  all  $i \neq j$  for some  $j \in \{1, \dots, n\}$ . For  $x_1^n = (x_1, \dots, x_n)$  and  $y_1^n = (y_1, \dots, y_n)$  differ only along coordinates  $j$ , then  $|G(x_1^n) - G(y_1^n)| \leq \frac{2B}{n}$ . So by the bounded difference inequality,

$$\mathbb{P}(|P_n - P|_{\mathcal{F}} \geq \mathbb{E}[|P_n - P|_{\mathcal{F}}] + t) \leq \exp\left(\frac{-nt^2}{2B^2}\right)$$

End of part 1 .

Step 2: Need to bound  $\mathbb{E}[|P_n - P|_{\mathcal{F}}]$ . We will handle this with a general result:

**Theorem 23.2** (*Symmetrization Inequalities*)

Let  $\mathcal{F}$  be a class of integrable functions on  $(\mathcal{X}, \mathcal{B})$ . Let  $\|\mathcal{R}_n\|_{\mathcal{F}} = \sup_{f \in \mathcal{F}} \frac{1}{n} |\sum_{i=1}^n \epsilon_i f(X_i)|$  where  $\tilde{X} = (X_1, \dots, X_n) \stackrel{iid}{\sim}$

$P \perp \tilde{\epsilon} = (\epsilon_1, \dots, \epsilon_n) \stackrel{iid}{\sim}$  Rademacher. Then for any convex non-decreasing function  $\phi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ ,

$$\begin{aligned} \mathbb{E}[\phi(|P_n - P|_{\mathcal{F}})] &\leq \mathbb{E}_{\tilde{X}, \tilde{\epsilon}}[\phi(2\|\mathcal{R}_n\|_{\mathcal{F}})] \\ \text{also } \mathbb{E}_{\tilde{X}, \tilde{\epsilon}}[\phi(\frac{1}{2}\|\mathcal{R}_n\|_{\tilde{\mathcal{F}}})] &\leq \mathbb{E}[\phi(|P_n - P|_{\mathcal{F}})] \\ \tilde{\mathcal{F}} &= \{f - \mathbb{E}[f(X_1)], f \in \mathcal{F}\} \end{aligned}$$

Returning to the proof, notice that  $\mathcal{R}_n(\mathcal{F}) = \mathbb{E}_{\tilde{X}, \tilde{\epsilon}}[\|\mathcal{R}_n\|_{\mathcal{F}}]$ . Take  $\phi(t) = t$  (meaning  $\phi$  is identity) to conclude that,

$$\mathbb{E}[|P_n - P|_{\mathcal{F}}] \leq 2\mathbb{E}_{\tilde{X}, \tilde{\epsilon}}[\|\mathcal{R}_n\|_{\mathcal{F}}] = 2\mathcal{R}_n(\mathcal{F})$$

This proves the theorem.

*Proof* of upper bound of symmetrization inequalities:

$$\mathbb{E}_{\tilde{X}}[\phi(|P_n - P|_{\mathcal{F}})] = \mathbb{E}_{\tilde{X}}[\phi(\sup_{f \in \mathcal{F}} \frac{1}{n} |\sum_{i=1}^n (f(X_i) - \mathbb{E}[f(X_i)])|)]$$

Let  $\tilde{Y} = (Y_1, \dots, Y_n) \stackrel{iid}{\sim} P$  be a ‘‘ghost’’ sample, where  $\tilde{Y} \perp \tilde{X}$  then  $\mathbb{E}[f(X_i)] = \mathbb{E}[f(Y_i)]$  so it

$$= \mathbb{E}_{\tilde{X}, \tilde{Y}}[\phi(\sup_{f \in \mathcal{F}} \frac{1}{n} |\sum_{i=1}^n (f(X_i) - \mathbb{E}_{Y_i}[f(Y_i)])|)]$$

By Jensen’s inequality:

$$\leq \mathbb{E}_{\tilde{X}, \tilde{Y}}[\phi(\sup_{f \in \mathcal{F}} \frac{1}{n} |\sum_{i=1}^n (f(X_i) - f(Y_i))|)].$$

Next  $f(X_i) - f(Y_i) \stackrel{d}{=} \epsilon_i(f(X_i) - f(Y_i))$  where  $\epsilon_i \perp X_i \& Y_i$  is Rademacher. Can now write

$$\leq \mathbb{E}_{\tilde{X}, \tilde{Y}, \tilde{\epsilon}}[\phi(\sup_{f \in \mathcal{F}} \frac{1}{n} |\sum_{i=1}^n \epsilon_i (f(X_i) - f(Y_i))|)]$$

Next use the triangle inequality and the fact that  $\phi$  is non-decreasing to get

$$\leq \mathbb{E}_{\tilde{X}, \tilde{Y}, \tilde{\epsilon}}[\phi(\sup_{f \in \mathcal{F}} \frac{1}{n} |\sum_{i=1}^n \epsilon_i f(X_i)| + \sup_{f \in \mathcal{F}} \frac{1}{n} |\sum_{i=1}^n \epsilon_i f(Y_i)|)]$$

Multiply and divide the two summands on the right-hand side by 2 (constructing a binary random variable),

$$\leq \mathbb{E}_{\tilde{X}, \tilde{Y}, \tilde{\epsilon}} \left[ \phi \left( \frac{1}{2} \sup_{f \in \mathcal{F}} \frac{2}{n} \left| \sum_i^n \epsilon_i f(X_i) \right| + \frac{1}{2} \sup_{f \in \mathcal{F}} \frac{2}{n} \left| \sum_i^n \epsilon_i f(Y_i) \right| \right) \right]$$

Then by Jensen's inequality,

$$\begin{aligned} &\leq \frac{1}{2} \mathbb{E}_{\tilde{X}, \tilde{Y}, \tilde{\epsilon}} \left[ \phi \left( \sup_{f \in \mathcal{F}} \frac{2}{n} \left| \sum_i^n \epsilon_i f(X_i) \right| \right) \right] + \frac{1}{2} \mathbb{E}_{\tilde{X}, \tilde{Y}, \tilde{\epsilon}} \left[ \phi \left( \sup_{f \in \mathcal{F}} \frac{2}{n} \left| \sum_i^n \epsilon_i f(Y_i) \right| \right) \right] \\ &= \mathbb{E}_{\tilde{X}, \tilde{\epsilon}} \left[ \phi \left( \sup_{f \in \mathcal{F}} \frac{2}{n} \left| \sum_i^n \epsilon_i f(X_i) \right| \right) \right] \\ &= \mathbb{E}_{\tilde{X}, \tilde{\epsilon}} \left[ \phi(2 \|\mathcal{R}_n\|_{\mathcal{F}}) \right] \end{aligned}$$

Our goal now (to use this result), is to upper bound,

$$\mathcal{R}_n(\mathcal{F}) = \mathbb{E}_{\tilde{X}, \tilde{\epsilon}} \left[ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_i^n \epsilon_i f(X_i) \right| \right]$$

## 23.2 VC Theory

**Definition 23.3** A class  $\mathcal{F}$  of functions on  $\mathcal{X}$  has polynomial discrimination with parameter  $\nu \geq 1$  if for each  $n \in \mathbb{N}$  and each  $n$ -tuple  $x_1^n = (x_1, \dots, x_n) \in \mathcal{X}$  then set  $\overline{\mathcal{F}}(x_i) = \{(f(x_1), \dots, f(x_n)) \in \mathbb{R}^n, f \in \mathcal{F}\} \subset \mathbb{R}^n$  has cardinality no larger than  $(n+1)^\nu$ .

This is a purely combinatorial property of class  $\mathcal{F}$ , nothing stochastic, always have  $n$ -tuple, then doing this for each function in the class, and get this bound.

**Lemma 23.4** Assume  $\mathcal{F}$  has polynomial discrimination with parameter  $\nu$ , then  $\forall n \in \mathbb{N}$  and for each  $x_1^n \in \mathcal{X}^n$ ,

$$\mathbb{E}_{\tilde{\epsilon}} \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \left| \sum_i^n \epsilon_i f(X_i) \right| \right] \leq 2D(x_1^n) \sqrt{\frac{\nu \log(n+1)}{n}}$$

where  $D(x_1^n) = \sup_{f \in \mathcal{F}} \sqrt{\frac{\sum_i^n f(x_i)^2}{n}}$ .

*Proof.* (Exercise) Use  $\sum \epsilon_i f(x_i) \sim SG(D(x_1^n)^2)$ .

**Corollary 23.5** If  $\mathcal{F}$  is a class uniformly bounded (by  $B > 0$ ), then  $\mathcal{R}_n(\mathcal{F}) \leq 2B \sqrt{\frac{\nu \log(n+1)}{n}}$

Example: Let  $\mathcal{F} = \{I\{(-\infty, z]\}, z \in \mathbb{R}\}$ ,

$$\|\mathcal{P}_n - P\|_{\mathcal{F}} = \sup_{z \in \mathbb{R}} |\hat{F}_n(z) - F_X(z)|$$

where  $\hat{F}_n(z)$  is the empirical CDF (this is the maximal difference between the CDF and empirical CDF). The  $\mathcal{F}$  has polynomial discrimination with parameter  $\nu = 1$ . Fix  $x_1^n \subset \mathbb{R}$ , let's order then

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$$

obtain  $n^{\text{th}}$  intervals

$$(-\infty, x_{(1)}], (x_{(1)}, x_{(2)}], \dots, (x_{(n)}, \infty)$$

as  $z$  varies over  $\mathbb{R}$  the function  $z \Rightarrow I\{(-\infty, z]\}$  will be 1 or 0 depending on how many interval  $z$  crosses in total, there are only  $(n + 1)$  possible realizations

$$\begin{aligned} \mathcal{F}(x_1^n) &\leq n + 1, \text{ all } x_1^n \\ \Rightarrow \mathbb{P}(\|P_n - P\|_{\mathcal{F}} \geq 4\sqrt{\frac{\log n}{n}} + t) &\leq \exp\frac{-nt^2}{2} \end{aligned}$$

Even stronger result DKW: (1990 Massart)

$$\mathbb{P}(\sup_{z \in \mathbb{R}} |\hat{F}_n(z) - F(z)| \geq t) \leq 2\exp\left(\frac{-nt^2}{2}\right)$$

This is a sharper result as the  $\log n$  term is gone, but this is only for the empirical CDF.