

Spring 2018: STA 6448  
Advanced Probability and Inference II  
Lecture 7

Yun Yang

- Uniform laws of large numbers

## Rademacher complexity

For any fixed collection  $x_1^n = (x_1, \dots, x_n)$  of points, consider the subset of  $\mathbb{R}^n$  given by

$$\mathcal{F}(x_1^n) = \left\{ (f(x_1), \dots, f(x_n)) \mid f \in \mathcal{F} \right\}.$$

Recall that the Rademacher complexity of this set (rescaled by  $n^{-1}$ ) is defined by

$$\mathcal{R}(\mathcal{F}(x_1^n)/n) = \mathbb{E}_\varepsilon \left[ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(x_i) \right| \right],$$

which is called the empirical Rademacher complexity.

### Definition

Given random samples  $X_1^n = (X_1, \dots, X_n)$ , the Rademacher complexity of the function class  $\mathcal{F}$  is defined as

$$\mathcal{R}_n(\mathcal{F}) = \mathbb{E}_X [\mathcal{R}(\mathcal{F}(x_1^n)/n)] = \mathbb{E}_{X, \varepsilon} \left[ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_i) \right| \right].$$

## A uniform law via Rademacher complexity

Rademacher complexity characterizes the typical largest correlation between a random noise vector and any function in the class  $\mathcal{F}$ , thereby the “complexity” of  $\mathcal{F}$ .

### Theorem

*Let  $\mathcal{F}$  be a class of functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  that is uniformly bounded by  $b > 0$ . Then for all  $n > 0$  and  $\delta \geq 0$ , we have*

$$\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \leq 2 \mathcal{R}_n(\mathcal{F}) + \delta$$

*with  $\mathbb{P}$  probability at least  $1 - 2 \exp\left(-\frac{n\delta^2}{8b^2}\right)$ . Consequently,  $\mathcal{R}_n(\mathcal{F}) = o(1)$  implies  $\mathcal{F}$  to be Glivenko-Cantelli.*

## Proof step one: Concentration around mean

Consider the function

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

It satisfies the bounded difference property: for all  $x_1, \dots, x_n, x'_k \in \mathbb{R}$ ,

$$\left| G(x_1, \dots, x_n) - G(x_1, \dots, x_{k-1}, x'_k, x_{k+1}, \dots, x_n) \right| \leq \frac{2\|f\|_\infty}{n} \leq \frac{2b}{n}.$$

Therefore, the bounded difference inequality implies the following holds with probability at least  $1 - 2 \exp \left( - \frac{nt^2}{8b^2} \right)$ ,

$$\left| \|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} - \mathbb{E}[\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}] \right| \leq t, \quad \text{for any } t > 0.$$

## Proof step two: Upper bound on mean

Applying the symmetrization technique.

Let  $(Y_1, \dots, Y_n)$  be a second independent copy of  $(X_1, \dots, X_n)$ .  
Then

$$\begin{aligned}\mathbb{E}[\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}] &= \mathbb{E}_X \left[ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \{f(X_i) - \mathbb{E}_{Y_i}[f(Y_i)]\} \right| \right] \\ &= \mathbb{E}_X \left[ \sup_{f \in \mathcal{F}} \left| \mathbb{E}_Y \left[ \frac{1}{n} \sum_{i=1}^n \{f(X_i) - f(Y_i)\} \right] \right| \right] \\ &\leq \mathbb{E}_{X,Y} \left[ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \{f(X_i) - f(Y_i)\} \right| \right],\end{aligned}$$

where the last step is due to Jensen's inequality.

## Proof step two: Upper bound on mean

Let  $\varepsilon_i$  be i.i.d. Rademacher random variables.

For any  $f \in \mathcal{F}$ , random variable  $\varepsilon_i(f(X_i) - f(Y_i))$  has the same distribution as  $f(X_i) - f(Y_i)$ . Consequently,

$$\begin{aligned}\mathbb{E}[\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}] &\leq \mathbb{E}_{X,Y} \left[ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \{f(X_i) - f(Y_i)\} \right| \right] \\ &\leq 2\mathbb{E}_{X,\varepsilon} \left[ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_i) \right| \right] = 2\mathcal{R}_n(\mathcal{F}).\end{aligned}$$

# Necessary conditions with Rademacher complexity

In the proof, we relate  $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$  with its symmetrized version

$$\|R_n\|_{\mathcal{F}} = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_i) \right|.$$

The stochastic process  $R_n$  over  $\mathcal{F}$  is known as the Rademacher process. What is lost in moving from  $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$  to  $\|R_n\|_{\mathcal{F}}$ ?  
Essentially nothing!

## Theorem

*For any convex non-decreasing function  $\Phi : \mathbb{R} \rightarrow \mathbb{R}$ , we have*

$$\mathbb{E}_{X, \varepsilon} \left[ \Phi \left( \frac{1}{2} \|R_n\|_{\bar{\mathcal{F}}} \right) \right] \leq \mathbb{E}[\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}] \leq \mathbb{E}_{X, \varepsilon} \left[ \Phi \left( 2 \|R_n\|_{\mathcal{F}} \right) \right],$$

*where  $\bar{\mathcal{F}} = \{f - \mathbb{E}[f] : f \in \mathcal{F}\}$  is the re-centered function class.*

# Necessary conditions with Rademacher complexity

## Corollary

*Let  $\mathcal{F}$  be a class of functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  that is uniformly bounded by  $b > 0$ . Then for all  $n > 0$  and  $\delta \geq 0$ , we have*

$$\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \geq \frac{1}{2} \mathcal{R}_n(\mathcal{F}) - \frac{b}{2\sqrt{n}} - \delta$$

*with  $\mathbb{P}$  probability at least  $1 - 2 \exp\left(-\frac{n\delta^2}{8b^2}\right)$ .*

Combined with the previous result, we obtain a two sided bound: with probability at least  $1 - 2 \exp\left(-\frac{n\delta^2}{8b^2}\right)$ ,

$$\frac{1}{2} \mathcal{R}_n(\mathcal{F}) - 2\delta \leq \|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \leq \mathcal{R}_n(\mathcal{F}) + \delta, \quad \text{for all } \delta > 0.$$



# Some upper bounds on Rademacher complexity

Given  $x_1^n = (x_1, \dots, x_n)$ , the size of

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

provides a sample-dependent measure of the complexity of  $\mathcal{F}$ .

## Definition

For a class  $\mathcal{F} \subset \{0, 1\}^{\mathcal{X}}$ , the growth function is

$$\Pi_{\mathcal{F}}(n) = \max \{|\mathcal{F}(x_1^n)| : x_1, \dots, x_n \in \mathcal{X}\}.$$

A class  $\mathcal{F}$  is said to have polynomial growth of order  $\nu \geq 1$  if

$$\Pi_{\mathcal{F}}(n) \leq (n + 1)^{\nu}, \quad \text{for all } n \geq 1.$$

# Controlling Rademacher complexity: growth function

## Theorem

For any  $x_1^n = (x_1, \dots, x_n)$ , let  $D(x_1^n) = \sup_{f \in \mathcal{F}} \sqrt{\frac{\sum_{i=1}^n f^2(x_i)}{n}}$  denote the  $\ell_2$  radius of  $\mathcal{F}(x_1^n)/\sqrt{n}$ . Then

$$\mathcal{R}(\mathcal{F}(x_1^n)/n) = \mathbb{E}_\varepsilon \left[ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(x_i) \right| \right] \leq D(x_1^n) \sqrt{\frac{2 \log(2 \Pi_{\mathcal{F}}(n))}{n}}.$$

In particular, if  $\mathcal{F}$  is uniformly bounded by  $b > 0$ , and has polynomial growth of order  $\nu \geq 1$ , then

$$\mathcal{R}_n(\mathcal{F}) \leq b \sqrt{\frac{2\nu \log(2(n+1))}{n}}.$$

Proof is left as a homework problem.

## Application: Classical Glivenko-Cantelli theorem

Recall the classical Glivenko-Cantelli theorem on the uniform convergence of CDFs:

$$\|\hat{F}_n - F\|_\infty \xrightarrow{\text{a.s.}} 0,$$

### Corollary

*Let  $F$  be the cdf and  $\hat{F}_n$  the empirical CDF, then*

$$\mathbb{P}\left[\|\hat{F}_n - F\|_\infty \geq \sqrt{\frac{2 \log(2(n+1))}{n}} + \delta\right] \leq 2e^{-\frac{n\delta^2}{8}} \quad \text{for all } \delta > 0,$$

*and hence  $\|\hat{F}_n - F\|_\infty \xrightarrow{\text{a.s.}} 0$ .*

*Proof:* Take  $\mathcal{F} = \{(-\infty, t] : t \in \mathbb{R}\}$ , then  $\mathcal{F}$  is uniformly bounded by 1, and has polynomial growth of order 1.

The bound is not tight (the  $\log(n+1)$  factor can be removed).

# Vapnik-Chervonenkis (VC) dimension

## Definition

A class  $\mathcal{F} \subset \{0, 1\}^{\mathcal{X}}$  shatters  $(x_1, \dots, x_d) \subset \mathcal{X}$  means  $|\mathcal{F}(x_1^d)| = 2^d$ .

The VC-dimension  $d_{VC}(\mathcal{F})$  is defined as the largest integer  $d$  for which there is some  $(x_1, \dots, x_d) \subset \mathcal{X}$  of  $d$  points that can be shattered by  $\mathcal{F}$ .

## Examples

- ▶  $\mathcal{F}_{\text{left}} = \{(-\infty, t] : t \in \mathbb{R}\}$  has VC-dim 1. It has polynomial growth of order 1.
- ▶  $\mathcal{F}_{\text{two}} = \{(s, t] : s, t \in \mathbb{R}\}$  has VC-dim 2. It has polynomial growth of order 2 (why?).

# Vapnik-Chervonenkis (VC) dimension

## Theorem (Sauer's Lemma)

*If  $d_{VC}(\mathcal{F}) \leq d$ , then*

$$\Pi_{\mathcal{F}}(n) \leq \sum_{k=1}^d \binom{n}{k} \leq (n+1)^d.$$

*Consequently, if  $d_{VC}(\mathcal{F}) < \infty$  (called VC class), then  $\mathcal{F}$  has polynomial growth of order  $d_{VC}(\mathcal{F})$ .*

*Proof:* See "Weak convergence and empirical processes: with applications to statistics", Section 2.6.1.

# Some useful results on Rademacher complexity

## Properties

1.  $\mathcal{F}_1 \subset \mathcal{F}_2$  implies  $\mathcal{R}_n(\mathcal{F}_1) \leq \mathcal{R}_n(\mathcal{F}_2)$ .
2. For any constant  $c \in \mathbb{R}$ ,  $\mathcal{R}_n(c\mathcal{F}) = |c| \mathcal{R}_n(\mathcal{F})$ .
3. For any fixed bounded function  $g$  (bounded by  $b$ ),  
 $|\mathcal{R}_n(\mathcal{F} + g) - \mathcal{R}_n(\mathcal{F})| \leq b \sqrt{2 \log 2/n}$ .
4.  $\mathcal{R}_n(\text{conv}(\mathcal{F})) = \mathcal{R}_n(\mathcal{F})$ , where  $\text{conv}(\mathcal{F})$  is the convex hull of  $\mathcal{F}$ .
5. If  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  is 1-Lipschitz continuous and satisfies  $\phi(0) = 0$ , then  $\mathcal{R}(\phi(\mathcal{F})) \leq 2\mathcal{R}(\mathcal{F})$ .

For a proof of the last claim, see “Probability in Banach Spaces” by Michel Ledoux and Michel Talagrand, Theorem 4.12.