

## Lecture 22: November 20

Lecturer: Alessandro Rinaldo

Scribe: Robin Dunn

**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.*

## 22.1 Sub-Gaussian Processes

We are often interested in bounding expressions of the form

$$\mathbb{E} \left[ \sup_{\theta \in \mathbb{T}} \theta^T \epsilon \right],$$

where  $\mathbb{T} \subset \mathbb{R}^n$ , and  $\epsilon = (\epsilon_1, \dots, \epsilon_n)$  is a vector of independent  $\text{SG}(\sigma^2)$  random variables.

**Example:** Suppose we have a class of the form  $\mathbb{T} = \mathcal{F}(x_1^n) = \{(f(x_1), \dots, f(x_n)) : f \in \mathcal{F}, x_i \in \mathbb{R}^d\}$  or  $\mathcal{A}(x_1^n) = \{A \cap x_1^n, A \in \mathcal{A}\}$ , where  $\mathcal{A}$  is a collection of subsets of  $\mathbb{R}^n$ .

- Suppose  $\epsilon = (\epsilon_1, \dots, \epsilon_n)$  is a vector of  $n$  independent Rademacher random variables. Then  $\mathcal{R}_n(\mathcal{F}) = \mathbb{E} \left[ \sup_{\theta \in \mathbb{T}} \theta^T \epsilon \right]$  is the Rademacher complexity of  $\mathbb{T}$  (or  $\mathcal{F}$ ).
- Suppose  $\epsilon = (\epsilon_1, \dots, \epsilon_n)$  is a vector of  $n$  independent  $N(0, 1)$  (or  $N(0, \sigma^2)$ ) random variables. Then  $\mathcal{G}_n(\mathcal{F}) = \mathbb{E} \left[ \sup_{\theta \in \mathbb{T}} \theta^T \epsilon \right]$  is the Gaussian complexity of  $\mathbb{T}$  (or  $\mathcal{F}$ ).

**Remark:** Rademacher and Gaussian complexities are sometimes of a similar order and sometimes of different orders.

If  $\mathbb{T} = \{\theta \in \mathbb{R}^d : \|\theta\|_2 \leq 1\}$ , then  $\mathcal{R}_n(\mathbb{T}) \approx \mathcal{G}_n(\mathbb{T}) \leq \sqrt{d}$ .

If  $\mathbb{T} = \{\theta \in \mathbb{R}^d : \|\theta\|_1 \leq 1\}$ , then  $\mathcal{R}_n(\mathbb{T}) = 1$  and  $\mathcal{G}_n(\mathbb{T}) \lesssim \sqrt{\log d}$ .

To show these results, we would use the facts  $\|x\|_2 = \sup_{v: \|v\|_2 \leq 1} v^T x = \sup_{v \in B_d(1)} v^T x$  and  $\|x\|_1 = \sup_{\|v\|_\infty \leq 1} v^T x$ .

Let  $\{X_\theta, \theta \in \mathbb{T}\}$  be a mean zero stochastic process indexed by  $\mathbb{T}$ . Similar to above, we may be interested in expression of the form  $\mathbb{E} \left[ \sup_{\theta \in \mathbb{T}} X_\theta \right]$ .

**Examples:**

- 1) Rademacher and Gaussian complexities. In the examples above, we could represent  $X_\theta = \epsilon^T \theta$ , where we are interested in  $\mathbb{E} \left[ \sup_{\theta \in \mathbb{T}} X_\theta \right]$ .
- 2) Non-parametric least-squares regression. We observe  $n$  pairs  $(Y_1, x_1), \dots, (Y_n, x_n)$  where  $x_1, \dots, x_n$  are deterministic points in  $[0, 1]$ . We assume that  $Y_i = f^*(x_i) + \epsilon_i$ , where  $(\epsilon_1, \dots, \epsilon_n) \stackrel{iid}{\sim} \text{SG}(\sigma^2)$  and

$f^* \in \mathcal{F}$  (a class of real-valued functions on  $[0, 1]$ ). Let  $\hat{f} \in \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum (Y_i - f(x_i))^2$  be the least squares estimator. We want to bound  $\mathbb{E} \left[ \underbrace{\frac{1}{n} \sum_{i=1}^n (\hat{f}(x_i) - f^*(x_i))^2}_{MSE} \right]$ .

Small  $MSE$  means that  $\hat{f}$  is a good approximation to  $f^*$ . To analyze the performance of  $\hat{f}$ , we start with the basic inequality:

$$\begin{aligned} MSE &\leq \frac{2}{n} \sum_{i=1}^n \epsilon_i (\hat{f}(x_i) - f^*(x_i)) \\ &\leq \frac{2}{\sqrt{n}} \sup_{f, g \in \mathcal{F}} |X_f - X_g| \end{aligned}$$

where  $X_f = \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f(x_i)$ . So

$$\mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (\hat{f}(x_i) - f^*(x_i))^2 \right] \leq \frac{2}{\sqrt{n}} \mathbb{E} \left[ \sup_{f, g \in \mathcal{F}} |X_f - X_g| \right].$$

Also, for any two functions  $f, g \in \mathcal{F}$ ,

$$\begin{aligned} V(X_f - X_g) &= \mathbb{E} \left[ \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i (f(x_i) - g(x_i)) \right)^2 \right] \\ &= \frac{1}{n} \mathbb{E} \left[ \left( \sum_{i=1}^n \epsilon_i (f(x_i) - g(x_i)) \right)^2 \right] \\ &\leq \frac{1}{n} \mathbb{E} \left[ \left( \sum_{i=1}^n \epsilon_i^2 \right) \left( \sum_{i=1}^n (f(x_i) - g(x_i))^2 \right) \right] \\ &= \frac{1}{n} \mathbb{E} \left[ \left( \sum_{i=1}^n \epsilon_i^2 \right) \sum_{i=1}^n (f(x_i) - g(x_i))^2 \right] \\ &\leq \frac{1}{n} \cdot n \sigma^2 \|f - g\|_2^2 \\ &= \sigma^2 \|f - g\|_2^2. \end{aligned}$$

- 3) Estimation in Wasserstein distance. Essentially, the Wasserstein distance is the amount of mass one must move from one distribution to another to make them equal.

Suppose  $P$  and  $Q$  are distributions on  $\mathbb{R}$ . The Wasserstein distance between  $P$  and  $Q$  is

$$W_1(P, Q) = \sup_{f \in \mathcal{F}} |Pf - Qf|,$$

where  $Pf = \mathbb{E}_{X \sim P}[f(X)]$  and  $\mathcal{F} = \{f : [0, 1] \rightarrow \mathbb{R}, f \text{ is 1-Lipschitz}\}$ . (That is, for any  $f \in \mathcal{F}$  and  $x, y \in [0, 1]$ ,  $|f(x) - f(y)| \leq |x - y|$ .)

An equivalent characterization is given by

$$W_1(P, Q) = \inf_{(x, y)} \mathbb{E}[|X - Y|, X \sim P, Y \sim Q].$$

We might want to use this metric to compare a true distribution to its empirical distribution. Suppose  $(X_1, \dots, X_n) \stackrel{iid}{\sim} P$  and  $P_n$  is the corresponding empirical measure. Then  $W_1(P_n, P) = \sup_{f \in \mathcal{F}} |X_f|$ , where  $X_f = P_n f - P f$ . So  $\mathbb{E}[X_f] = 0$  for all  $f$ . Then we see

$$\mathbb{E}[W_1(P_n, P)] = \mathbb{E}\left[\sup_{f \in \mathcal{F}} |X_f|\right].$$

### 22.1.1 Sub-Gaussian Processes

**Definition 22.1** (*Sub-Gaussian process.*) A zero-mean stochastic process  $\{X_\theta : \theta \in \mathbb{T}\}$  is sub-Gaussian with respect to metric  $d$  on  $\mathbb{T}$  if for  $\theta, \theta' \in \mathbb{T}$  and  $\lambda \in \mathbb{R}$ ,

$$\mathbb{E}\left[e^{\lambda(X_\theta - X_{\theta'})}\right] \leq \exp\left[\frac{\lambda^2}{2} d^2(\theta, \theta')\right].$$

Equivalently, for  $\theta, \theta' \in \mathbb{T}$ ,  $X_\theta - X_{\theta'} \in \text{SG}(d^2(\theta, \theta'))$ .  $d(\cdot)$  is called the canonical metric. In the case of Gaussian random variables, the canonical metric is given by  $d(\theta, \theta') = \sqrt{V(X_\theta - X_{\theta'})}$ .

By Hoeffding's inequality for sub-Gaussians,

$$\mathbb{P}(|X_\theta - X_{\theta'}| \geq t) \leq 2 \exp\left\{-\frac{t^2}{2d^2(\theta, \theta')}\right\}.$$

**Examples:**

- 1) Rademacher and Gaussian complexities. In these cases,  $\mathbb{T} \subseteq \mathbb{R}^n$ . These processes are sub-Gaussian with respect to  $d(\theta, \theta') = \|\theta - \theta'\|_2$  on  $\mathbb{T}$  because

$$V(X_\theta - X_{\theta'}) = V(\epsilon^T \theta - \epsilon^T \theta') \leq \|\theta - \theta'\|_2^2 \sigma^2. \quad (22.1)$$

(This proves that  $X_\theta - X_{\theta'} \in \text{SG}(\|\theta - \theta'\|_2^2 \sigma^2)$ .) In the case where  $\epsilon_i$  is Rademacher,  $\sigma^2 = 1$ . In the case where  $\epsilon_i \sim N(0, \sigma^2)$ ,  $\sigma^2$  in equation 22.1 is the same  $\sigma^2$  as the normal variance.

- 2) Non-parametric least squares regression. As before, we define  $X_f = \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f(x_i)$ , where  $x_1, \dots, x_n$  are deterministic.  $X_f$  is SG, and so is  $X_f - X_g$ . Previously, we showed  $V(X_f - X_g) \leq \sigma^2 \|f - g\|_2^2$ . In this case, we can use the canonical distance  $d(f, g) = \sqrt{\frac{1}{n} \sum_{i=1}^n (f(x_i) - g(x_i))^2}$ .
- 3) Wasserstein distance. In this problem, we had  $X_f = P_n f - P f$ . This is an SG process with respect to  $d(f, g) = \frac{\|f - g\|_\infty}{\sqrt{n}}$ . This is an exercise, and the result can be obtained by using Azuma-Hoeffding.

### 22.1.2 Metric Entropy

**Definition 22.2** (*Metric entropy.*) Let  $\mathbb{T} \subseteq \mathbb{R}^n$  and let  $d$  be a distance metric on  $\mathbb{T}$ . For  $\delta > 0$ , the metric entropy of  $\mathbb{T}$  with respect to  $d$  is given by  $\log \mathcal{N}(\mathbb{T}, \delta)$ , where  $\mathcal{N}(\mathbb{T}, \delta)$  is the  $\delta$ -covering number of  $\mathbb{T}$ .

**Definition 22.3** (*Diameter of  $\mathbb{T}$ .*) Let  $\mathbb{T} \subseteq \mathbb{R}^n$  and let  $d$  be a distance metric on  $\mathbb{T}$ . The diameter of the set  $\mathbb{T}$  is given by  $D = \sup_{\theta, \theta' \in \mathbb{T}} d(\theta, \theta')$ .

**Proposition 22.4** (1-step discretization bound.) Assume  $\{X_\theta : \theta \in \mathbb{T}\}$  is a SG process with respect to  $d$ . Then for all  $\delta \in (0, D]$ ,

$$\mathbb{E} \left[ \sup_{\theta, \theta' \in \mathbb{T}} |X_\theta - X_{\theta'}| \right] \leq 2 \mathbb{E} \left[ \sup_{\substack{\gamma, \gamma' \in \mathbb{T} \\ d(\gamma, \gamma') \leq \delta}} |X_\gamma - X_{\gamma'}| \right] + 4D \sqrt{\log \mathcal{N}(\mathbb{T}, \delta)}.$$

$\delta$  is a tuning parameter. As  $\delta$  decreases, the first term decreases and the second term increases.

Remarks:

- 1) For arbitrary  $\theta_0 \in \mathbb{T}$ ,  $\mathbb{E} \left[ \sup_{\theta \in \mathbb{T}} X_\theta \right] = \mathbb{E} \left[ \sup_{\theta \in \mathbb{T}} (X_\theta - X_{\theta_0}) \right] \leq \mathbb{E} \left[ \sup_{\theta, \theta' \in \mathbb{T}} (X_\theta - X_{\theta'}) \right]$ .
- 2) Constants are not optimal.

**Proof:** Let  $\theta_1, \dots, \theta_N$  be a minimal  $\delta$ -cover of  $\mathbb{T}$ , where  $N = \mathcal{N}(\mathbb{T}, \delta)$ . Then for all  $\theta \in \mathbb{T}$ , there exists  $j$  ( $1 \leq j \leq N$ ) such that  $d(\theta, \theta_j) \leq \delta$ .

Fix  $\theta \in \mathbb{T}$ . Choose  $j$  such that  $d(\theta, \theta_j) \leq \delta$ . Then

$$\begin{aligned} X_\theta - X_{\theta_1} &= X_\theta - X_{\theta_j} + X_{\theta_j} - X_{\theta_1} \\ &\leq \sup_{\substack{\gamma, \gamma' \in \mathbb{T} \\ d(\gamma, \gamma') \leq \delta}} (X_\gamma - X_{\gamma'}) + \max_i |X_{\theta_i} - X_{\theta_1}|. \end{aligned}$$

We can obtain a similar bound for  $X_{\theta_1} - X_{\theta'}$ , where  $\theta'$  is another point in  $\mathbb{T}$ .

Adding up and using the fact that  $\theta$  and  $\theta'$  are arbitrary,

$$\sup_{\theta, \theta' \in \mathbb{T}} (X_\theta - X_{\theta'}) \leq 2 \sup_{\substack{\gamma, \gamma' \in \mathbb{T} \\ d(\gamma, \gamma') \leq \delta}} (X_\gamma - X_{\gamma'}) + 2 \max_i |X_{\theta_i} - X_{\theta_1}|.$$

To finish the proof, we will take the expectation of both sides. Since  $X_{\theta_i} - X_{\theta_1} \in \text{SG}(D^2)$ , we know  $\mathbb{E} [\max_i |X_{\theta_i} - X_{\theta_1}|] \leq 2D \sqrt{\log \mathcal{N}(\mathbb{T}, \delta)}$ . (See the maximal inequality from the 9-13 lecture notes.) So

$$\begin{aligned} \mathbb{E} \left[ \sup_{\theta, \theta' \in \mathbb{T}} |X_\theta - X_{\theta'}| \right] &\leq 2 \mathbb{E} \left[ \sup_{\substack{\gamma, \gamma' \in \mathbb{T} \\ d(\gamma, \gamma') \leq \delta}} (X_\gamma - X_{\gamma'}) \right] + 2 \mathbb{E} \left[ \max_i |X_{\theta_i} - X_{\theta_1}| \right] \\ &\leq 2 \mathbb{E} \left[ \sup_{\substack{\gamma, \gamma' \in \mathbb{T} \\ d(\gamma, \gamma') \leq \delta}} |X_\gamma - X_{\gamma'}| \right] + 4D \sqrt{\log \mathcal{N}(\mathbb{T}, \delta)}. \end{aligned}$$

■

**Applications:** For  $\mathbb{T} \subseteq \mathbb{R}^n$  and  $\delta \in (0, D]$  (where  $D$  is the diameter of  $\mathbb{T}$ ), define

$$\tilde{\mathbb{T}}(\delta) := \{\gamma - \gamma' | \gamma, \gamma' \in \mathbb{T}, \|\gamma - \gamma'\|_2 \leq \delta\}.$$

Then where  $\epsilon$  is a vector of Rademacher random variables and  $d(\cdot)$  is Euclidean distance,

$$\mathcal{R}_n(\tilde{\mathbb{T}}(\delta)) = \mathbb{E} \left[ \sup_{\substack{\gamma, \gamma' \in \mathbb{T} \\ d(\gamma, \gamma') \leq \delta}} \epsilon^T (\gamma - \gamma') \right] \leq \mathbb{E} [\|\epsilon\|_2 \delta] \leq \sqrt{n} \delta.$$

The same inequality holds for  $\mathcal{G}_n(\widetilde{\mathbb{T}}(\delta))$  if  $\epsilon$  is a vector of  $N(0, 1)$  random variables.

Applying the 1-step discretization bound, we see

$$\mathbb{E} \left[ \sup_{\theta, \theta' \in \mathbb{T}} \epsilon^T (\theta - \theta') \right] \lesssim \min_{\delta \in (0, D]} \{ \delta \sqrt{n} + \sqrt{\log \mathcal{N}(\mathbb{T}, \delta)} \}$$

(up to constants). Again, as  $\delta \rightarrow 0$ ,  $\delta \sqrt{n} \rightarrow 0$  and  $\sqrt{\log \mathcal{N}(\mathbb{T}, \delta)}$  increases (often to infinity). To balance, we set  $\delta \sqrt{n} = \sqrt{\log \mathcal{N}(\mathbb{T}, \delta)}$  and solve for  $\delta$ .

## References

- [W17] M. WAINWRIGHT, “High-dimensional statistics: A Non-asymptotic Viewpoint. (Draft),” Chapter 5. 2017.