## SYS 6016/4582: Machine LearningSpring 2017Lecture 5Instructor: Quanquan GuDate: Feb 2<sup>nd</sup>, 2017Scriber: Robert Carroll

Last time we introduced **Bounded Difference Inequality** to bound the supremum of the difference between the empirical risk and expected risk, i.e.,

$$\sup_{h \in \mathcal{H}} \left| \hat{R}_n(h) - R(h) \right| \le \max \left\{ \sup_{h \in \mathcal{H}} \hat{R}_n(h) - R(h), \sup_{h \in \mathcal{H}} R(h) - \hat{R}_n(h) \right\}$$

With the fact that

$$\mathbb{P}\left(\sup_{h\in\mathcal{H}}\left|\hat{R}_{n}(h)-R(h)\right|\geq t\right)\leq\mathbb{P}\left(\underbrace{\sup_{h\in\mathcal{H}}\hat{R}_{n}(h)-R(h)}_{G_{n}(\mathcal{H})}\geq t\right)+\mathbb{P}\left(\underbrace{\sup_{h\in\mathcal{H}}R(h)-\hat{R}_{n}(h)}_{G_{n}'(\mathcal{H})}\geq t\right),$$

we can bound the term  $\sup_{h \in \mathcal{H}} |\hat{R}_n(h) - R(h)|$  by bounding  $G_n(\mathcal{H})$  and  $G'_n(\mathcal{H})$  respectively. Since  $G_n(\mathcal{H})$  and  $G'_n(\mathcal{H})$  can basically be bounded by the same term, we only show how to bound  $G_n(\mathcal{H})$ .

We restate the Bounded Difference Inequality here.

## **Theorem 1 (Bounded Difference Inequality)** Let f be a function satisfying

$$|f(z_1, z_2, ..., z_i, ..., z_n) - f(z_1, z_2, ..., z'_i, ..., z_n)| \le c_i$$

for all i = 1, 2, ..., n and for all  $z_1, z_2, ..., z_i, z'_i, ..., z_n$ . Let  $Z_1, Z_2, ..., Z_n$  be independent random variables. We have that

$$\mathbb{P}\left(f\left(Z_{1}, Z_{2}, ..., Z_{n}\right) - \mathbb{E}[f(Z_{1}, Z_{2}, ..., Z_{n})]\right) \ge t\right) \le 2\exp\left(\frac{-2t^{2}}{\sum_{i=1}^{n} c_{i}^{2}}\right).$$

## Roadmap to bound $G_n(\mathcal{H})$

Step 1: Prove  $G_n(\mathcal{H}) \leq \mathbb{E}G_n(\mathcal{H}) + \sqrt{\frac{\log(1/\delta)}{2n}}$  with probability at least  $1 - \delta$ . Step 2: Bound  $\mathbb{E}G_n(\mathcal{H})$  by symmetrization.

Step 1: Bounding  $G_n(\mathcal{H})$ 

**Theorem 2** Let  $G_n(\mathcal{H}) = \sup_{h \in \mathcal{H}} \hat{R}_n(h) - R(h)$ , we have the following holds with probability at least  $1 - \delta$ 

$$G_n(\mathcal{H}) - \mathbb{E}[G_n(\mathcal{H})] \le \sqrt{\frac{\log(1/\delta)}{2n}}.$$

**Proof:** Let  $Z_i = (X_i, Y_i) \stackrel{i.i.d.}{\sim} P_{X,Y}$  and  $f(Z_1, Z_2, ..., Z_n) = G_n(\mathcal{H})$ . First we verify that the function  $f(Z_1, ..., Z_n)$  satisfies the conditions in Theorem 1. For any *i*, the following equality holds:

$$\left| f\left(Z_{1}, Z_{2}, ..., Z_{i}, ..., Z_{n}\right) - f\left(Z_{1}, Z_{2}, ..., Z_{i}', ..., Z_{n}\right) \right|$$
  
=  $\left| \sup_{h \in \mathcal{H}} \left( \hat{R}_{n}(h) - R(h) \right) - \sup_{h \in \mathcal{H}} \left( \hat{R}_{n}'(h) - R(h) \right) \right|,$ (1)

where  $\hat{R}'_{n}(h)$  is the empirical risk when the observation  $Z_{i}$  is replaced by  $Z'_{i}$ , i.e.,

$$\hat{R}'_n(h) = \frac{1}{n} \left[ \sum_{j \neq i}^n \mathbb{1} \left\{ h(X_j) \neq Y_j \right\} + \mathbb{1} \left\{ h(X'_i) \neq Y'_i \right\} \right].$$

Thus we can obtain

$$\left|\sup_{h\in\mathcal{H}} \left(\hat{R}_n(h) - R(h)\right) - \sup_{h\in\mathcal{H}} \left(\hat{R}'_n(h) - R(h)\right)\right| \le \sup_{h\in\mathcal{H}} \left|\hat{R}_n(h) - \hat{R}'_n(h)\right|.$$
 (2)

Using the definition of the empirical risk, the R.H.S. of (2) can be rewritten as:

$$\begin{split} \sup_{h \in \mathcal{H}} \left| \hat{R}_n(h) - \hat{R}'_n(h) \right| \\ &= \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \mathbbm{1} \left\{ h(X_i) \neq Y_i \right\} - \frac{1}{n} \left[ \sum_{j \neq i}^n \mathbbm{1} \left\{ h(X_j) \neq Y_j \right\} + \mathbbm{1} \left\{ h(X'_i) \neq Y'_i \right\} \right] \right| \\ &= \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \mathbbm{1} \left\{ h(X_i) \neq Y_i \right\} - \frac{1}{n} \mathbbm{1} \left\{ h(X'_i) \neq Y'_i \right\} \right| = \frac{1}{n} \triangleq c_i, \end{split}$$

where in the third equality we use the fact that indicator functions are either 0 or 1, so the maximal difference between any two indicator functions is 1. Hence we proved that  $G_n(\mathcal{H})$  satisfies the conditions in 1 with  $c_i = 1/n$ . Applying Theorem 1 gives

$$\mathbb{P}(G_n(\mathcal{H}) - \mathbb{E}G_n(\mathcal{H}) \ge t) \le \exp\left(\frac{-2t^2}{\sum_{i=1}^n c_i^2}\right)$$
$$= \exp\left(\frac{-2t^2}{\sum_{i=1}^n (1/n)^2}\right)$$
$$= \exp\left(-2nt^2\right).$$

Let  $\delta = \exp(-2nt^2)$ . Solving for t gives

$$t = \sqrt{\frac{\log(1/\delta)}{2n}}$$

We then have with probability at least  $1 - \delta$ 

$$G_n(\mathcal{H}) - \mathbb{E}G_n(\mathcal{H}) \le \sqrt{\frac{\log(1/\delta)}{2n}},$$

equivalently

$$G_n(\mathcal{H}) \leq \mathbb{E}G_n(\mathcal{H}) + \sqrt{\frac{\log(1/\delta)}{2n}}.$$

This completes the proof.  $\blacksquare$ 

## Step 2: Bounding $\mathbb{E}G_n(\mathcal{H})$

Before presenting the result we first lay out some definitions and lemmas that are essential in the proof of the main theorem.

**Definition 1 (Symmetric Random Variable)** Z is a symmetric random variable if  $\mathbb{P}(Z \ge t) = \mathbb{P}(-Z \ge t) = \mathbb{P}(-Z \le -t)$ . In other words, multiplying by -1 doesn't change the distribution.

The following lemma shows how to construct a symmetric random variable from i.i.d. random variables.

**Lemma 1** If X, Y are two independent and identically distributed random variables, then Z = X - Y is a symmetric random variable.

The following lemma is essential to bound  $\mathbb{E}G_n(\mathcal{H})$ , which is called the symmetrization technique.

**Lemma 2** If  $Z_1, Z_2, \ldots, Z_n$  are symmetric random variables, and  $\sigma_1, \sigma_2, \ldots, \sigma_n$  are Rademacher random variables, *i.e.*,

$$\sigma_i = \begin{cases} 1, & \text{with probability } 1/2, \\ -1, & \text{with probability } 1/2. \end{cases}$$

In addition, suppose  $\sigma_1, \ldots, \sigma_n$  are independent of  $Z_1, \ldots, Z_n$ , then  $\sum_{i=1}^n Z_i$  has the same distribution as  $\sum_{i=1}^n \sigma_i Z_i$ .

Based on the definitions and lemmas listed above, the bound of  $\mathbb{E}G_n(\mathcal{H})$  can be represented as the following theorem.

**Theorem 3** We have

$$\mathbb{E}G_n(\mathcal{H}) \leq 2\mathbb{E}_{X,Y}\mathbb{E}_{\sigma} \sup_{h \in \mathcal{H}} \left[\frac{1}{n} \sum_{i=1}^n \sigma_i \mathbb{1}\left\{h(X_i) \neq Y_i\right\}\right].$$

**Proof:** We prove using symmetrization technique with a "ghost" sample  $(\widetilde{X}_i, \widetilde{Y}_i)_{i=1}^n \stackrel{i.i.d}{\sim} P_{X,Y}$  that is independent of the training sample

First, we lay out several claims, which are useful to bound  $\mathbb{E}G_n(\mathcal{H})$ .

**Claim 1** Given the "ghost" sample  $(\widetilde{X}_i, \widetilde{Y}_i)_{i=1}^n$ , we have the following hold

$$R(h) = \mathbb{E}_{\widetilde{X},\widetilde{Y}}\left[\frac{1}{n}\sum_{i=1}^{n}\mathbb{1}\left\{h(\widetilde{X}_{i})\neq\widetilde{Y}_{i}\right\}\right].$$

By the definition of  $G_n(\mathcal{H})$ , we have

$$\mathbb{E}G_{n}(\mathcal{H}) = \mathbb{E}_{X,Y} \left[ \sup_{h \in \mathcal{H}} \left( \hat{R}_{n}(h) - R(h) \right) \right] \\ = \mathbb{E}_{X,Y} \left[ \sup_{h \in \mathcal{H}} \left( \hat{R}_{n}(h) - \mathbb{E}_{\widetilde{X},\widetilde{Y}} \left[ \tilde{R}_{n}(h) \right] \right) \right],$$
(3)

where

$$\tilde{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\left\{h(\widetilde{X}_i) \neq \widetilde{Y}_i\right\} \text{ and } (\widetilde{X}_i, \widetilde{Y}_i), \stackrel{i.i.d}{\sim} P_{X,Y}.$$

In addition we have  $(\widetilde{X}_i, \widetilde{Y}_i)_{i=1}^n$  are independent from  $(X_i, Y_i)_{i=1}^n$ , where  $(\widetilde{X}_i, \widetilde{Y}_i)_{i=1}^n$  is called the "ghost" sample.

According to (3), we have

$$\begin{split} & \mathbb{E}_{X,Y} \left[ \sup_{h \in \mathcal{H}} \left( \hat{R}_n(h) - \mathbb{E}_{\widetilde{X},\widetilde{Y}} \left[ \tilde{R}_n(h) \right] \right) \right] \\ &= \mathbb{E}_{X,Y} \left[ \sup_{h \in \mathcal{H}} \mathbb{E}_{\widetilde{X},\widetilde{Y}} \left( \hat{R}_n(h) - \tilde{R}_n(h) \right) \right] \\ &\leq \mathbb{E}_{X,Y} \left[ \mathbb{E}_{\widetilde{X},\widetilde{Y}} \sup_{h \in \mathcal{H}} \left( \hat{R}_n(h) - \tilde{R}_n(h) \right) \right] \\ &= \mathbb{E}_{X,Y} \mathbb{E}_{\widetilde{X},\widetilde{Y}} \left[ \sup_{h \in \mathcal{H}} \left( \frac{1}{n} \sum_{i=1}^n \mathbbm{1} \left\{ h(X_i) \neq Y_i \right\} - \frac{1}{n} \sum_{i=1}^n \mathbbm{1} \left\{ h(\widetilde{X}_i) \neq \widetilde{Y}_i \right\} \right) \right], \end{split}$$

where the second inequality comes from Jensen's Inequality and the fact that the supremum function is convex. This implies that

$$\mathbb{E}G_{n}(\mathcal{H}) \leq \mathbb{E}_{X,Y}\mathbb{E}_{\widetilde{X},\widetilde{Y}}\left[\sup_{h\in\mathcal{H}}\left(\frac{1}{n}\sum_{i=1}^{n}\mathbb{1}\left\{h(X_{i})\neq Y_{i}\right\}-\frac{1}{n}\sum_{i=1}^{n}\mathbb{1}\left\{h(\widetilde{X}_{i})\neq \widetilde{Y}_{i}\right\}\right)\right]$$
$$=\mathbb{E}_{X,Y}\mathbb{E}_{\widetilde{X},\widetilde{Y}}\left[\sup_{h\in\mathcal{H}}\left(\frac{1}{n}\sum_{i=1}^{n}\left[\mathbb{1}\left\{h(X_{i})\neq Y_{i}\right\}-\mathbb{1}\left\{h(\widetilde{X}_{i})\neq \widetilde{Y}_{i}\right\}\right]\right)\right].$$
(4)

Since  $\mathbb{1}\left\{h(X_i) \neq Y_i\right\}$ , and  $\mathbb{1}\left\{h(\widetilde{X}_i) \neq \widetilde{Y}_i\right\}$  are two independent and identically distributed random variables, according to Lemma 1, we have  $\mathbb{1}\left\{h(X_i) \neq Y_i\right\} - \mathbb{1}\left\{h(\widetilde{X}_i) \neq \widetilde{Y}_i\right\}$ 

is symmetric. According to Lemma 2, (4) can be further upper bounded by

$$\mathbb{E}G_{n}(\mathcal{H}) \leq \mathbb{E}_{X,Y}\mathbb{E}_{\widetilde{X},\widetilde{Y}}\left[\sup_{h\in\mathcal{H}}\left(\frac{1}{n}\sum_{i=1}^{n}\left[\mathbb{1}\left\{h(X_{i})\neq Y_{i}\right\}-\mathbb{1}\left\{h(\widetilde{X}_{i})\neq\widetilde{Y}_{i}\right\}\right]\right)\right].$$

$$=\mathbb{E}_{X,Y}\mathbb{E}_{\widetilde{X},\widetilde{Y}}\mathbb{E}_{\sigma}\left[\sup_{h\in\mathcal{H}}\left(\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}\left[\mathbb{1}\left\{h(X_{i})\neq Y_{i}\right\}-\mathbb{1}\left\{h(\widetilde{X}_{i})\neq\widetilde{Y}_{i}\right\}\right]\right)\right]$$

$$\leq \mathbb{E}_{X,Y}\mathbb{E}_{\widetilde{X},\widetilde{Y}}\mathbb{E}_{\sigma}\left[\sup_{h\in\mathcal{H}}\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}\mathbb{1}\left\{h(X_{i})\neq Y_{i}\right\}+\sup_{h\in\mathcal{H}}\frac{1}{n}\sum_{i=1}^{n}(-\sigma_{i})\mathbb{1}\left\{h(\widetilde{X}_{i})\neq\widetilde{Y}_{i}\right\}\right],$$

where the last line comes from the inequality that  $\sup_x (f(x)+g(x)) \leq \sup_x f(x)+\sup_x g(x)$ . By the linearity of expectation we have

$$\mathbb{E}G_{n}(\mathcal{H}) \leq \mathbb{E}_{X,Y} \mathbb{E}_{\widetilde{X},\widetilde{Y}} \mathbb{E}_{\sigma} \bigg[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \mathbb{1} \Big\{ h(X_{i}) \neq Y_{i} \Big\} \bigg] \\ + \mathbb{E}_{X,Y} \mathbb{E}_{\widetilde{X},\widetilde{Y}} \mathbb{E}_{\sigma} \bigg[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} (-\sigma_{i}) \mathbb{1} \Big\{ h(\widetilde{X}_{i}) \neq \widetilde{Y}_{i} \Big\} \bigg].$$

Since  $\sigma_i$ 's are symmetric random variables,  $-\sigma_i$  has the same distribution with  $\sigma_i$ . Then changing  $-\sigma_i$  into  $\sigma_i$  in the second term does not change the value of the whole expression. Hence we have Furthermore, we can get

$$\mathbb{E}G_{n}(\mathcal{H}) \leq \mathbb{E}_{X,Y}\mathbb{E}_{\widetilde{X},\widetilde{Y}}\mathbb{E}_{\sigma}\left[\sup_{h\in\mathcal{H}}\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}\mathbb{1}\left\{h(X_{i})\neq Y_{i}\right\}\right] \\ + \mathbb{E}_{X,Y}\mathbb{E}_{\widetilde{X},\widetilde{Y}}\mathbb{E}_{\sigma}\left[\sup_{h\in\mathcal{H}}\frac{1}{n}\sum_{i=1}^{n}(-\sigma_{i})\mathbb{1}\left\{h(\widetilde{X}_{i})\neq\widetilde{Y}_{i}\right\}\right] \\ = \mathbb{E}_{\sigma}\left[\mathbb{E}_{X,Y}\sup_{h\in\mathcal{H}}\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}\mathbb{1}\left\{h(X_{i})\neq Y_{i}\right\} + \mathbb{E}_{\widetilde{X},\widetilde{Y}}\sup_{h\in\mathcal{H}}\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}\mathbb{1}\left\{h(\widetilde{X}_{i})\neq\widetilde{Y}_{i}\right\}\right] \\ = 2\mathbb{E}_{X,Y}\mathbb{E}_{\sigma}\sup_{h\in\mathcal{H}}\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}\mathbb{1}\left\{h(X_{i})\neq Y_{i}\right\}.$$
(5)

This completes the proof.  $\blacksquare$ 

Next, we introduce the concept of Rademacher Complexity, which measures the complexity of a class of functions and has a wide range of applications in machine learning theory. Note that the following definition is tailored to 0-1 loss function.

**Definition 2 (Rademacher Complexity)** For a given function class  $\mathcal{F}$ , its Rademacher complexity is defined as

$$\mathcal{R}_n(\mathcal{F}) = \mathbb{E}_Z \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(Z_i) \right],$$

where  $\sigma_i$ 's are *i.i.d.* Rademacher random variables, and  $Z_i$ 's are *i.i.d.* random variables.

Based on the definition of Rademacher Complexity, we can also define the Empirical Rademacher Complexity.

**Definition 3 (Empirical Rademacher Complexity)** For a given hypothesis class  $\mathcal{F}$  and a set of observations  $Z_i$ 's, the empirical Rademacher complexity is

$$\hat{\mathcal{R}}_n(\mathcal{F}) = \mathbb{E}_{\sigma} \bigg[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(Z_i) \ \Big| Z_{1:n} \bigg].$$

where  $\sigma_i$ 's are *i.i.d.* Rademacher random variables, and  $Z_{1:n}$  is a shorthand notation for  $Z_1, Z_2, ..., Z_n$ .