

### 3 - Symmetrization

Last time:

- Generalization by
  - 1 - concentration
  - 2 - discretization + union bound
  - 3 - uniform convergence  $\Rightarrow$  generalization

Today: A new technique to replace step 2, namely symmetrization.

- Slight modification uniform conv  $\Rightarrow$  generalization

\* Before  $\mathbb{P}(R(\hat{f}) - R(f_*) \geq \epsilon) \leq \mathbb{P}\left(\sup_{\mathcal{F}} |\hat{R}(f) - R(f)| \geq \frac{\epsilon}{2}\right)$

which was due to

$$\begin{aligned}\epsilon &\leq R(\hat{f}) - R(f_*) \leq R(f) - \hat{R}(f) + 0 + \hat{R}(f_*) - R(f_*) \\ &\leq \sup_{\mathcal{F}} R(f) - \hat{R}(f) + \sup_{\mathcal{F}} \hat{R}(f) - R(f) \\ &\leq \frac{\epsilon}{2} \leq \frac{\epsilon}{2}\end{aligned}$$

$$\Rightarrow \mathbb{P}(R(\hat{f}) - R(f_*) \geq \epsilon) \leq \mathbb{P}\left(\left\{\sup_{\mathcal{F}} R(f) - \hat{R}(f) \geq \frac{\epsilon}{2}\right\} \cup \left\{\sup_{\mathcal{F}} \hat{R}(f) - R(f) \geq \frac{\epsilon}{2}\right\}\right)$$

$$\begin{aligned}(\text{by union bound}) \quad &\leq \mathbb{P}\left(\sup_{\mathcal{F}} R(f) - \hat{R}(f) \geq \frac{\epsilon}{2}\right) \\ &+ \mathbb{P}\left(\sup_{\mathcal{F}} \hat{R}(f) - R(f) \geq \frac{\epsilon}{2}\right)\end{aligned}$$

$$\begin{aligned}(\text{by symmetry}) \quad &= 2 \cdot \mathbb{P}\left(\sup_{\mathcal{F}} \underbrace{\hat{R}(f) - R(f)}_{\text{empirical process}} \geq \frac{\epsilon}{2}\right)\end{aligned}$$

$\Rightarrow$  Need to bound RHS.

**Theorem** (Generalization by RC): For a fnc class  $\mathcal{F}$ , define  $\mathcal{G} = \{ (x,y) \rightarrow \ell((x,y), f) : f \in \mathcal{F} \}$ . If the loss fnc satisfies  $0 \leq \ell \leq B$ , then with prob at least  $1-\delta$ , we have

$$R(\hat{f}) - R(f^*) \leq 4 R(\mathcal{G}) + B \sqrt{\frac{2 \log 2/\delta}{n}}$$

- Remarks:**
- RC is a complexity measure of a fnc class.
  - Rate depends on  $R(\mathcal{G})$ .
    - $\mathcal{G}_f \in \mathcal{G}$  depends on  $f \in \mathcal{F}$ . We expect  $R(\mathcal{G}) \approx R(\mathcal{F})$ ?
    - We hope, as  $n \uparrow$   $R(\mathcal{G}) \downarrow$ .

**proof:**

- Strategy:
- 1 - Concentration (~~Hoeffding~~, McDiarmid's)
  - 2 - ~~Discretization + union bound~~ Symmetrisation
  - 3 - Unif conv.  $\Rightarrow$  generalization

Goal: Bound the empirical process.

### Step 1: Concentration

**Theorem** (McDiarmid's Inequality): Let  $g$  be a function satisfying the "bounded difference" property,

$$\forall j \in [n] \quad |g(x_1, \dots, x_j, \dots, x_n) - g(x_1, \dots, x'_j, \dots, x_n)| \leq c_j.$$

Then, for  $z_1, z_2, \dots, z_n$  independent r.v.'s

$$\mathbb{P}\left(g(z_1, \dots, z_n) - \mathbb{E}g(z_1, \dots, z_n) \geq \epsilon\right) \leq \exp\left\{-\frac{2\epsilon^2}{\sum_{j=1}^n c_j^2}\right\}.$$

Example (Hoeffding's inequality):  $z_1, \dots, z_n$  independent and  
 $\forall i \quad a_i \leq z_i \leq b_i$ .

$$g(z_1, \dots, z_n) = \frac{1}{n} \sum_{i=1}^n z_i$$

$$\begin{aligned} |g(z_1, \dots, z_j, \dots, z_n) - g(z_1, \dots, z'_j, \dots, z_n)| &= \left| \frac{1}{n} \sum_{i=1}^n z_i - \frac{1}{n} \sum_{i=1}^n z_i - \frac{1}{n} (z'_j - z_j) \right| \\ &= \frac{1}{n} |z_j - z'_j| \leq \frac{b_j - a_j}{n} := c_j \end{aligned}$$

By McDiarmid's inequality

$$\mathbb{P} \left( \frac{1}{n} \sum_{i=1}^n z_i - \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n z_i \right] \geq \epsilon \right) \leq \exp \left\{ - \frac{2n^2 \epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2} \right\}$$

Goal: Bound  $\sup_f \hat{R}(f) - R(f)$

$$g(\underline{z}_1, \dots, \underline{z}_n) = \sup_f \hat{R}(f) - R(f)$$

$$= \sup_f \frac{1}{n} \sum_{i=1}^n \ell(\underbrace{(x_i, y_i)}_{\underline{z}_i}, f) - \mathbb{E} \ell(\underbrace{(x, y)}_{\underline{z}}, f)$$

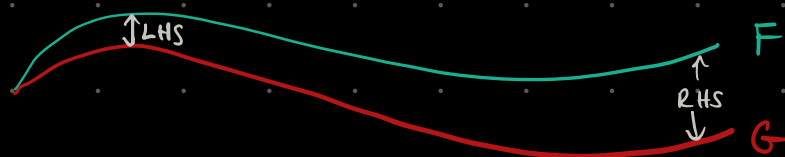
$$= \sup_f \frac{1}{n} \sum_{i=1}^n \ell(z_i, f) - \mathbb{E} \ell(z, f)$$

$$|g(\dots z_j \dots) - g(\dots z'_j \dots)|$$

$$* = \left| \sup_f \frac{1}{n} \sum_{i=1}^n \ell(z_i, f) - \mathbb{E} \ell(z, f) \right|$$

$$= \left\{ \sup_f \frac{1}{n} \sum_{i=1}^n \ell(z_i, f) - \mathbb{E} \ell(z, f) - \frac{\ell(z_j, f) - \ell(z'_j, f)}{n} \right\}$$

Fact:  $\left| \sup_x F(x) - \sup_x G(x) \right| \leq \sup_x |F(x) - G(x)|$



$$* \leq \sup_f \left| \frac{\ell(z_j, f) - \ell(z'_j, f)}{n} \right| \leq \frac{B}{n} := c_f.$$

(since  $0 \leq \ell \leq B$ )

By McDiarmid's Inequality:

$$\mathbb{P} \left( \underbrace{\sup_f \hat{R}(f) - R(f)}_{\text{what we want}} - \underbrace{\mathbb{E} \left[ \sup_f \hat{R}(f) - R(f) \right]}_{\text{what we want}} \geq t \right) \leq \exp \left\{ - \frac{2nt^2}{B^2} \right\}$$

Need to show it is small.

Step 2: Symmetrization (To bound  $\rightarrow$ )

- Basic idea:  $X$  is a r.v. and  $X'$  is its iid copy.

$$\Rightarrow X \stackrel{d}{=} X' \quad \text{let } g \text{ be any fnc.}$$

$$\Rightarrow g(X) - g(X') \stackrel{d}{=} g(X') - g(X)$$

$$\stackrel{d}{=} -1 (g(X) - g(X'))$$

$$\stackrel{d}{=} \sigma (g(X) - g(X'))$$

where  $\sigma$  is a Rademacher r.v. :  $\mathbb{P}(\sigma = +1) = \frac{1}{2}$   
 $\mathbb{P}(\sigma = -1) = \frac{1}{2}$

- In our case, data  $\mathcal{D} = \{ (x_1, y_1), \dots, (x_n, y_n) \}$  is r.v.  
 $= \{ z_1, \dots, z_n \}$

- Introduce iid copy of the dataset (ghost data)

$$\mathcal{D}' = \{z_1', \dots, z_n'\}$$

$z_i$ 's and  $z_i'$ 's are iid.

- Now, we have 2 empirical risks, 1 population risk.

$$i) \hat{R}(f; \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n \ell(z_i, f)$$

$$ii) \hat{R}(f; \mathcal{D}') = \frac{1}{n} \sum_{i=1}^n \ell(z_i', f)$$

$$\begin{aligned} \mathbb{E} \hat{R}(f; \mathcal{D}) &= \mathbb{E} \hat{R}(f; \mathcal{D}') \\ &= R(f) \end{aligned}$$

- Notice that  $R(f) = \mathbb{E}[\hat{R}(f; \mathcal{D})] = \mathbb{E}[\hat{R}(f; \mathcal{D}') | \mathcal{D}]$ .

Goal: Bound  $\mathbb{E} \left[ \sup_f \hat{R}(f) - R(f) \right]$ .

$$\mathbb{E} \left[ \sup_f \hat{R}(f; \mathcal{D}) - R(f) \right] = \mathbb{E} \left[ \sup_f \left\{ \hat{R}(f; \mathcal{D}) - \mathbb{E}[\hat{R}(f; \mathcal{D}')] \right\} \right]$$

$$= \mathbb{E} \left[ \sup_f \left\{ \hat{R}(f; \mathcal{D}) - \mathbb{E}[\hat{R}(f; \mathcal{D}') | \mathcal{D}] \right\} \right]$$

$$= \mathbb{E} \left[ \sup_f \mathbb{E} \left[ \hat{R}(f; \mathcal{D}) - \hat{R}(f; \mathcal{D}') | \mathcal{D} \right] \right]$$

$$\left\{ \text{by } \sup \mathbb{E} \leq \mathbb{E} \sup \right\} \leq \mathbb{E} \left[ \mathbb{E} \left[ \sup_f \hat{R}(f; \mathcal{D}) - \hat{R}(f; \mathcal{D}') | \mathcal{D} \right] \right]$$

$$\left\{ \begin{array}{l} \text{by tower property} \\ \mathbb{E}[\mathbb{E}[X|Y]] = \mathbb{E}[X] \end{array} \right\} = \mathbb{E} \left[ \sup_f \hat{R}(f; \mathcal{D}) - \hat{R}(f; \mathcal{D}') \right]$$

$$= \mathbb{E} \left[ \sup_f \frac{1}{n} \sum_{i=1}^n \underbrace{\ell(z_i, f) - \ell(z_i', f)}_{\triangleq \sigma_i \{ \ell(z_i, f) - \ell(z_i', f) \}} \right]$$

$$= \mathbb{E} \left[ \sup_{\mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \{ \ell(z_i, f) - \ell(z_i', f) \} \right]$$

$$\left\{ \begin{array}{l} \text{Fact: } \sup_z \{ F(z) + G(z) \} \\ \leq \sup_z F(z) + \sup_z G(z) \end{array} \right\}$$

$$\leq \mathbb{E} \left[ \sup_{\mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \ell(z_i, f) \right] + \mathbb{E} \left[ \sup_{\mathcal{F}} \frac{1}{n} \sum_{i=1}^n -\sigma_i \ell(z_i, f) \right]$$

$$\sigma_i \stackrel{d}{=} -\sigma_i$$

$$= 2 \mathbb{E} \left[ \sup_{\mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \ell(z_i, f) \right] = **$$

Definition (Rademacher Complexity): For a fnc class

$\mathcal{F} = \{ f: \mathcal{Z} \rightarrow \mathbb{R} \}$  and a dataset  $\mathcal{D} = \{ z_1, \dots, z_n \}$ ,

\* RC is defined as

$$\mathcal{R}(\mathcal{F}) = \mathbb{E} \left[ \sup_{\mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \right] \text{ where}$$

$\sigma_i$ 's are iid Rademacher r.v.'s.

\* Empirical RC is defined as

$$\hat{\mathcal{R}}(\mathcal{F}) = \mathbb{E} \left[ \sup_{\mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \mid z_{1:n} \right]$$

$$** = 2 \mathbb{E} \left[ \sup_{\mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \ell(z_i, f) \right] \leq 2 \mathcal{R}(\mathcal{F}) \quad \text{no!}$$

$$\leq 2 \mathcal{R}(\mathcal{G})$$

$$\text{where } \mathcal{G} = \{ z \rightarrow \ell(z, f) : f \in \mathcal{F} \}.$$

Step 3: Uniform convergence  $\Rightarrow$  generalization

$$\mathbb{P}(R(\hat{f}) - R(f_*) \geq \epsilon) \leq \underline{2} \cdot \mathbb{P}\left(\sup_{\mathcal{F}} \hat{R}(f) - R(f) \geq \underline{\frac{\epsilon}{2}}\right)$$

- By Step 1:  $\ast \mathbb{P}\left(\sup_{\mathcal{F}} \hat{R}(f) - R(f) \geq \mathbb{E}\left[\sup_{\mathcal{F}} \hat{R}(f) - R(f)\right] + t\right) \leq e^{-\frac{2nt^2}{B^2}}$

- By Step 2:  $\ast \mathbb{E}\left[\sup_{\mathcal{F}} \hat{R}(f) - R(f)\right] \leq 2 \cdot \mathcal{R}(g)$ .

- By Steps 1 and 2:

$$2 \mathbb{P}\left(\sup_{\mathcal{F}} \hat{R}(f) - R(f) \geq \underbrace{t + 2\mathcal{R}(g)}_{\frac{\epsilon}{2}}\right) \leq 2e^{-\frac{2nt^2}{B^2}} = \delta$$

$$\Rightarrow \delta = 2e^{-\frac{2nt^2}{B^2}} \Rightarrow t = B \sqrt{\frac{\log 2/\delta}{2n}}$$

$$\Rightarrow \frac{\epsilon}{2} := t + 2\mathcal{R}(g) \Rightarrow \epsilon = 4\mathcal{R}(g) + B \sqrt{\frac{2 \log 2/\delta}{n}} \quad \square$$