# 2 - Uniform Convergence

Today: A non-trivial setting where we use uniform conv.

- Recall the bound on excess risk:

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

excess risk

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

- We'll bound $\mathbb{P}\left(\underbrace{\sup_{\mathcal{F}} \hat{R}(f) - R(f)}_{\text{empirical process}} \geq \frac{\epsilon}{2}\right)$ which will

imply a bound on the first term by symmetry.

**Theorem** (Rademacher Complexity): Define $\mathcal{G} = \left\{(x,y) \rightarrow \ell((x,y),f) : f \in \mathcal{F}\right\}$.
If loss satisfies $0 \leq \ell \leq 1$, then with probability at least $1-\delta$,

$$R(\hat{f}) - R(f_*) \leq 2 \mathcal{R}(\mathcal{G}) + \sqrt{\frac{2 \log \frac{2}{\delta}}{n}}.$$

Rademacher Complexity (RC)

— RC is a complexity measure of a fnc class.
— Rate depends on $R(\mathcal{G})$.
  — $g_f \in \mathcal{G}$ depends on $f \in \mathcal{F}$. We expect $R(\mathcal{G}) \simeq R(\mathcal{F})$?
  — We hope, as $n \uparrow$ $R(\mathcal{G}) \downarrow$.

**proof:**

Strategy:  1 – Concentration ( ~~Hoeffding~~ , Mc Diarmid's )

   2 – ~~union bound~~ Symmetrisation

   3 – Unif conv. $\Rightarrow$ generalisation

Goal : Bound the empirical process.

### Step 1 : Concentration

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

※  $\forall j \in [n]$  $|g(x_1,\ldots,x_j,\ldots,x_n) - g(x_1,\ldots,x_j',\ldots,x_n)| \leq c_j$.

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

$$\mathbb{P}\left( g(z_1,\ldots,z_n) - \mathbb{E}g(z_1,\ldots,z_n) \geq \epsilon \right) \leq \exp\left\{ -\frac{2\epsilon^2}{\sum_{j=1}^{n} c_j^2} \right\}.$$
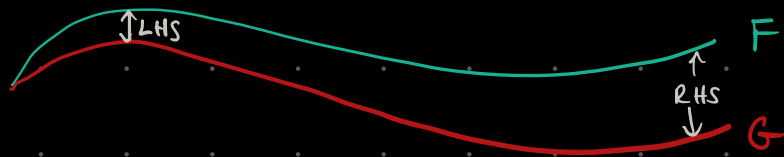
Application: ( Hoeffding's Inequality )

Goal: Bound $\sup_{\mathcal{F}} \hat{R}(f) - R(f)$

$$g(z_1, \ldots, z_n) = \sup_{\mathcal{F}} \hat{R}(f) - R(f)$$

$$= \sup_{\mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} \ell(\underbrace{(x_i, y_i), f}_{z_i}) - \mathbb{E}\, \ell(\underbrace{(x,y), f}_{z})$$

$$\left| g(\ldots z_j \ldots) - g(\ldots z_j' \ldots) \right|$$

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

$$\left. - \sup_{\mathcal{F}} \left\{ \frac{1}{n} \sum_{i=1}^{n} \ell(z_i, f) - \mathbb{E}\, \ell(z, f) - \frac{\ell(z_j, f) - \ell(z_j', f)}{n} \right\} \right|$$

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



↑LHS    F    ↑RHS    G

$$* \le \sup_{\mathcal{F}} \left| \frac{\ell(z_j, f) - \ell(z_j', f)}{n} \right| \le \frac{1}{n} := c_j.$$

$$\left( \text{Since} \quad 0 \le \ell \le 1 \right)$$

By McDiarmid's Inequality:

$$\mathbb{P}\left( \sup_{\mathcal{F}} \hat{R}(f) - R(f) - \mathbb{E}\left[ \sup_{\mathcal{F}} \hat{R}(f) - R(f) \right] \ge t \right) \le \exp\left\{ -2nt^2 \right\}$$

what we want

→ Need to show it is small.

Step 1: Symmetrization ( To bound ⌐⌐ )

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

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

$$\Rightarrow g(x) - g(x') \overset{d}{=} g(x') - g(x)$$

$$\overset{d}{=} -1 \, (g(x) - g(x'))$$

$$\overset{d}{=} \sigma \, (g(x) - g(x'))$$

where $\sigma$ is a Rademacher r.v.:

$$\mathbb{P}(\sigma = +1) = \frac{1}{2}$$
$$\mathbb{P}(\sigma = -1) = \frac{1}{2}$$

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

$$= \{ z_1, \dots, z_n \}$$

- Introduce iid copy of the dataset $D' = \{ z_1', \dots, z_n' \}$

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

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

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

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

$\mathbb{E}\hat{R}(f; D) = \mathbb{E}\hat{R}(f, D')$

$= R(f)$

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

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

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

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

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

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

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

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

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

$\left\{\begin{array}{l}\text{Fact: } \sup_x \{F(x) + G(x)\}\\ \quad\leq \sup_x F(x) + \sup_x G(x)\end{array}\right\}$ $\leq \mathbb{E}\left[\sup_f \frac{1}{n}\sum_{i=1}^{n}\sigma_i \ell(z_i,f)\right] + \mathbb{E}\left[\sup_f \frac{1}{n}\sum_{i=1}^{n}-\sigma_i \ell(z_i',f)\right]$

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

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

**Definition** (Radenacher Complexity): For a fnc class
$\mathcal{F} = \{ f : Z \longrightarrow \mathbb{R} \}$ and a dataset $D = \{ z_1 \cdots z_n \}$

* RC is defined as

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

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

* Empirical RC is defined as

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

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

$$= 2\,R(\mathcal{G})$$

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

**Step 3:** Uniform convergence $\Rightarrow$ generalization

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

– By Step 1: $\quad * \; \mathbb{P}\left( \sup_{\mathcal{F}} \widehat{R}(f) - R(f) \geq \mathbb{E}[\sup_{\mathcal{F}} \widehat{R}(f) - R(f)] + t \right) \leq e^{-2nt^2}$

– By Step 2: $\quad * \; \mathbb{E}\left[ \sup_{\mathcal{F}} \widehat{R}(f) - R(f) \right] \leq 2 \cdot R(\mathcal{G})$

- By Steps 1 and 2:

$$2\,\mathbb{P}\left(\sup_{\mathcal{F}} \hat{R}(f) - R(f) \geqslant t + 2R(\mathcal{G})\right) \leq 2\,e^{-2nt^2}$$

$$\underbrace{\phantom{t + 2R(\mathcal{G})}}_{\varepsilon/2} \qquad := \delta$$

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

$$\Rightarrow \quad \frac{\varepsilon}{2} := t + 2R(\mathcal{G}) \quad \Rightarrow \quad \varepsilon = 4R(\mathcal{G}) + \sqrt{\frac{2\log 2/\delta}{n}} \qquad \boxed{}$$

— **Generalization via RC**

Goal: i) Relate $R(\mathcal{G})$ to $R(\mathcal{F})$.

ii) $R(\mathcal{F})$ decays w/ $n$.

(i) — Theorem (Talagrand's Contraction Principle): Let $g$ be a $L$-Lipschitz cont. func. and $g \circ \mathcal{F} = \{g \circ f : f \in \mathcal{F}\}$, then

$$R(g \circ \mathcal{F}) \leq L \cdot R(\mathcal{F}).$$

— RC of Constrained Linear Models

(ii) — Goal: $R(\mathcal{F}) = O(1/\sqrt{n})$.

Theorem (RC of Linear Models): Let $\mathcal{F} = \{f(x) = \langle x, \theta \rangle : \|\theta\| \leq r\}$. Then, i) $\hat{R}(\mathcal{F}) \leq \frac{r}{n}\sqrt{\sum_{i=1}^{n}\|x_i\|^2}$

ii) If $\mathbb{E}[\|x_i\|^2] \leq \kappa^2$, then $R(\mathcal{F}) \leq \frac{r \cdot \kappa}{\sqrt{n}}$

Remarks: 1- If we combine this bound w/ previous examples, we achieve generalization.

2 - Notice $\kappa = O(\sqrt{d})$ so $\mathcal{R}(\mathcal{F}) \leq r\sqrt{\frac{d}{n}}$.

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

$$= \mathbb{E}\left[\sup_{\|\theta\|\leq r} \frac{1}{n}\sum_{i=1}^{n}\sigma_i \langle x_i, \theta\rangle \mid x_{1:n}\right]$$

$$= \mathbb{E}\left[\sup_{\|\theta\|\leq r} \left\langle \frac{1}{n}\sum_{i=1}^{n}\sigma_i x_i, \theta\right\rangle \mid x_{1:n}\right]$$

$\boxed{\sup_{\|\theta\|\leq r} \langle\theta, v\rangle = r\cdot\|v\|}$

$$= r\,\mathbb{E}\left[\left\| \frac{1}{n}\sum_{i=1}^{n}\sigma_i x_i\right\| \mid x_{1:n}\right]$$

(by Jensen's Ineq.)

$$\leq \frac{r}{n}\,\mathbb{E}\left[\left\|\sum_{i=1}^{n}\sigma_i x_i\right\|^2 \mid x_{1:n}\right]^{\frac{1}{2}}$$

$$= \frac{r}{n}\,\mathbb{E}\left[\sum_{i=1}^{n}\|x_i\|^2 + \sum_{i\neq j}\sigma_i\sigma_j\langle x_i, x_j\rangle \mid x_{1:n}\right]^{\frac{1}{2}}$$

$$= \frac{r}{n}\,\sqrt{\sum_{i=1}^{n}\|x_i\|^2}$$

ii) $\mathcal{R}(\mathcal{F}) = \mathbb{E}\hat{\mathcal{R}}(\mathcal{F}) \leq \frac{r}{n}\,\mathbb{E}\sqrt{\sum_{i=1}^{n}\|x_i\|^2}$

(by Jensen's Ineq)

$$\leq \frac{r}{n}\sqrt{\mathbb{E}\sum_{i=1}^{n}\|x_i\|^2}$$

$$\leq \frac{r}{n}\sqrt{\sum_{i=1}^{n}\mathbb{E}\{\|x_i\|^2\}}$$

$$\leq \frac{r\cdot\kappa}{\sqrt{n}}.$$