

# 1 - Uniform Convergence $\Rightarrow$ Generalization

- Supervised learning:  $(y, x) \sim p(y, x)$  iid pairs.  
 $y \in \mathcal{Y} = \mathbb{R}$ ,  $x \in \mathcal{X} = \mathbb{R}^d$

- Observe data:  $(y_i, x_i) \sim p$  for  $i = 1, 2, \dots, n$ .

- **Goal:** Find a function  $f \in \mathcal{F}$  s.t.  $f: \mathcal{X} \rightarrow \mathcal{Y}$   
 $y \approx f(x)$ .

\* Need to choose  $\mathcal{F} = \{f\}$ , the function class  
and the loss func  $l: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$

- **Goal+:** Find  $f^* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} R(f) = \mathbb{E} \left[ l(y, x), f \right]$   
↑  
over  $p$

\* Empirical Risk Minimization (ERM)

- Observe data  $\mathcal{D} = \{(y_i, x_i) : i = 1, \dots, n\}$ , then

$\hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \hat{R}(f) := \frac{1}{n} \sum_{i=1}^n l(y_i, x_i), f$   
↓ estimates

are these close?  $\rightarrow f^* \leftarrow \underset{f \in \mathcal{F}}{\operatorname{argmin}} R(f) = \mathbb{E} \left[ l(y, x), f \right]$

Ex (MLE):  $\mathcal{F} = \{f_\theta : \theta \in \Theta\}$ ,  $l = -\log p_\theta(y|x)$

-  $\hat{\theta} = \underset{\Theta}{\operatorname{argmin}} \hat{R}(\theta) = \frac{1}{n} \sum_{i=1}^n -\log p_\theta(y_i|x_i)$

-  $\theta_* = \underset{\Theta}{\operatorname{argmin}} R(\theta) = \mathbb{E} \left[ -\log p_\theta(y|x) \right]$

$\hat{R}(\hat{\theta})$ : training error  
 $R(\hat{\theta})$ : test error  
 Def (Excess risk):  $R(\hat{\theta}) - R(\theta^*)$

Generalization = small excess risk

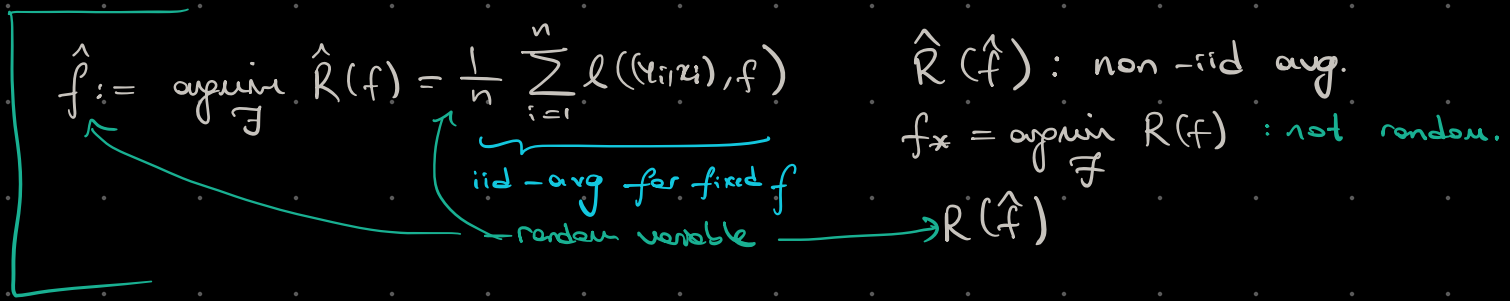
### Uniform Convergence

Goal: Understand generalization (= small excess risk), in a non-asymptotic sense.

$$\mathbb{P} \left( \underbrace{R(\hat{f}) - R(f^*)}_{\text{excess risk}} > \epsilon \right) \leq \delta$$

bad event:
↓  
small prob

Def (Uniform conv.):  $\sup_{f \in \mathcal{F}} |\hat{R}(f) - R(f)| \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 0$



$$\begin{aligned} \epsilon \leq R(\hat{f}) - R(f^*) &= \underbrace{\left\{ R(\hat{f}) - \hat{R}(\hat{f}) \right\}}_{\substack{\text{hard to handle} \\ \text{since } \hat{f} \text{ is random} \\ \Rightarrow \text{non-iid avg.}}} + \underbrace{\left\{ \hat{R}(\hat{f}) - \hat{R}(f^*) \right\}}_{\substack{\leq 0 \\ \text{since } \hat{f} = \operatorname{argmin}_{\mathcal{F}} \hat{R}(f)}}} + \underbrace{\left\{ \hat{R}(f^*) - R(f^*) \right\}}_{\substack{\text{iid avg} \\ \Rightarrow \text{ez to handle}}} \\ &\leq \sup_{f \in \mathcal{F}} R(f) - \hat{R}(f) + 0 + \sup_{f \in \mathcal{F}} \hat{R}(f) - R(f) \\ &\geq \epsilon/2 \qquad \qquad \qquad \geq \epsilon/2 \end{aligned}$$

$$\left[ \leq 2 \cdot \sup_{f \in \mathcal{F}} |\hat{R}(f) - R(f)| \quad \text{if } \xrightarrow{P} 0 \text{ then generalization.} \right]$$

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

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

$\Rightarrow$  Need to bound RHS.

\* It usually suffices to bound only one.  
(by symmetry)

- Generalization for Finite Function Classes. ( $|\mathcal{F}| < \infty$ )  
(warm-up)

**Theorem:** If  $|\mathcal{F}| < \infty$  and  $\delta \in [0, 1]$ , then

$$\mathbb{P}\left(R(\hat{f}) - R(f_*) \geq \sqrt{\frac{2}{n} (\log 2|\mathcal{F}| + \log \frac{1}{\delta})}\right) \leq \delta$$

$\swarrow$  sample size
 $\downarrow$  complexity of  $\mathcal{F}$ 
 $\searrow$  confidence level

Remark: - Generalization error rate:  $O\left(\sqrt{\frac{\log |\mathcal{F}|}{n}}\right)$

Proof: Strategy:

- 1- Concentration
- 2- Union bound
- 3- Generalization

→ Lemma (Hoeffding's Ineq.): Let  $z_1, z_2, \dots, z_n$  be indep. r.v.'s such that  $a_i \leq z_i \leq b_i$  almost surely.

Then,  $\forall \epsilon > 0$  for  $S_n = \frac{1}{n} \sum_{i=1}^n z_i$

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

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

Note that  $\hat{R}(f) = \frac{1}{n} \sum_{i=1}^n \ell((x_i, z_i), f)$  is like

$$S_n = \frac{1}{n} \sum_{i=1}^n z_i \quad \begin{array}{c} 0 \leq z \leq 1 \\ \downarrow \quad \downarrow \\ a_i \quad b_i \end{array}$$

1- Concentration: Fix  $f \in \mathcal{F}$ , then

$$\begin{aligned} \text{(by Hoeffding)} \quad \mathbb{P}(\hat{R}(f) - R(f) \geq \frac{\epsilon}{2}) &\leq \exp \left\{ - \frac{2n^2 \epsilon^2}{\sum_{i=1}^n 1} \right\} \\ &= \exp \left\{ - \frac{2n \epsilon^2}{4} \right\} \end{aligned}$$

$$= \exp \left\{ - \frac{n \epsilon^2}{2} \right\}$$

2- Uniform Convergence (via union bound):

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}} \hat{R}(f) - R(f) \geq \frac{\varepsilon}{2}\right) = \mathbb{P}\left(\bigcup_{f \in \mathcal{F}} \left\{ \hat{R}(f) - R(f) \geq \frac{\varepsilon}{2} \right\}\right)$$

$$\text{(by union bound)} \leq \sum_{f \in \mathcal{F}} \mathbb{P}\left(\hat{R}(f) - R(f) \geq \frac{\varepsilon}{2}\right)$$

$$\text{(by Hoeffding)} \leq |\mathcal{F}| e^{-n\varepsilon^2/2}$$

3- Generalization

$$\begin{aligned} \mathbb{P}\left(R(\hat{f}) - R(f_*) \geq \varepsilon\right) &\leq \mathbb{P}\left(\sup_{\mathcal{F}} R(f) - \hat{R}(f) \geq \frac{\varepsilon}{2}\right) \\ &\quad + \mathbb{P}\left(\sup_{\mathcal{F}} \hat{R}(f) - R(f) \geq \frac{\varepsilon}{2}\right) \\ &\leq 2 \cdot |\mathcal{F}| \exp\left\{-\frac{n\varepsilon^2}{2}\right\} := \delta \end{aligned}$$

$$\Rightarrow \log \frac{\delta}{2|\mathcal{F}|} = -\frac{n\varepsilon^2}{2} \Rightarrow \varepsilon = \sqrt{\frac{2}{n} \log \frac{2|\mathcal{F}|}{\delta}} \quad \square$$

- Remarks:
1. Means: with prob  $1-\delta$ ,  $R(\hat{f}) - R(f_*) \leq O\left(\sqrt{\frac{\log|\mathcal{F}| + \log\delta^{-1}}{n}}\right)$
  2. Choose  $\delta = O(|\mathcal{F}|^{-1})$   
Conv. rate becomes  $O\left(\sqrt{\frac{\log|\mathcal{F}|}{n}}\right)$
  3. Covers bounded loss, cannot cover square loss.
  4. Bound fails when  $|\mathcal{F}| = \infty$ !