

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.

* _____ the loss $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$

$$-\text{Goal+}: \text{Find } f_* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} R(f) = \mathbb{E}[\ell(y, f(x))]$$

\uparrow
 over $(y, x) \sim P$
 cannot!

* Empirical Risk Minimization (ERM)

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

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

hope ss \downarrow estimator
 $f_* \leftarrow \underset{\theta \in \Theta}{\operatorname{argmin}} R(\theta) = \mathbb{E}[\ell(y, f_\theta(x))]$

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

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

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

$$\left. \begin{array}{l} \hat{R}(\hat{\theta}) : \text{training error} \\ R(\hat{\theta}) : \text{test error} \end{array} \right\} \quad \text{Def (Excess risk)}: R(\hat{\theta}) - R(\theta_*)$$

Generalization = small excess risk

- Uniform Convergence

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

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

↓
small prob

bad event.

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

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

iid-avg for fixed f
random variable

$\hat{R}(\hat{f})$: non-iid avg.
 $f^* = \underset{\mathcal{F}}{\operatorname{argmin}} R(f)$: not random.

$$R(\hat{f}) - R(f^*) = \{R(\hat{f}) - \hat{R}(\hat{f})\} + \{\hat{R}(\hat{f}) - \hat{R}(f^*)\} + \{\hat{R}(f^*) - R(f^*)\}$$

hard to handle

(\hat{f} is random
- non-iid avg -)

≤ 0

$$\hat{f} = \underset{\mathcal{F}}{\operatorname{argmin}} \hat{R}(f)$$

ez to handle
as f^* is fixed.
- iid avg -

$$\leq \sup_{f \in \mathcal{F}} |R(f) - \hat{R}(f)| + 0 + \sup_{f \in \mathcal{F}} |\hat{R}(f) - R(f)|$$

$$\leq 2 \cdot \sup_{f \in \mathcal{F}} |\hat{R}(f) - R(f)|$$

If $\xrightarrow{P} 0$ (uniform convergence)

we have small excess risk (generalization)

$$\{ R(\hat{f}) - R(f_*) \geq \epsilon \} \stackrel{\text{A}}{\subseteq} \left\{ \sup_{f \in \mathcal{F}} | \hat{R}(f) - R(f) | \geq \frac{\epsilon}{2} \right\} \stackrel{\text{B}}{\subseteq}$$

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

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

Theorem: If $|\mathcal{F}| < \infty$ and $\ell \in [0, B]$, then

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

↓
 sample size ↓
 complexity of \mathcal{F} ↓
 confidence level

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

Proof:

→ 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 $S_n = \frac{1}{n} \sum_{i=1}^n z_i$

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

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

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

$$S_n = \frac{1}{n} \sum_{i=1}^n z_i$$

↓
 a_i ↓
 b_i

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

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

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

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

Remark: looks exponential in n , in fact it implies $\frac{1}{m}$ rate.

2- Uniform Convergence (via union bound)

$$\begin{aligned} \mathbb{P}\left(\sup_{f \in \mathcal{F}} |\hat{R}(f) - R(f)| \geq \frac{\epsilon}{2}\right) &= \mathbb{P}\left(\bigcup_{f \in \mathcal{F}} \left\{|\hat{R}(f) - R(f)| \geq \frac{\epsilon}{2}\right\}\right) \\ &\stackrel{\text{by union bound}}{\leq} \sum_{f \in \mathcal{F}} \mathbb{P}\left(|\hat{R}(f) - R(f)| \geq \frac{\epsilon}{2}\right) \\ &\quad \underbrace{\leq 2 \exp\left\{-\frac{n\epsilon^2}{2B^2}\right\}}_{\leq 2|\mathcal{F}| \exp\left\{-\frac{n\epsilon^2}{2B^2}\right\}} \\ &\stackrel{\text{by Hoeffding}}{\leq} 2|\mathcal{F}| \exp\left\{-\frac{n\epsilon^2}{2B^2}\right\} \end{aligned}$$

3 - Generalization

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

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

- Remarks:
1. Means: with prob $1-\delta$, $R(\hat{f}) - R(f_*) \leq B \sqrt{\frac{2}{n} \log \frac{2|\mathcal{F}|}{\delta}}$
 $= O\left(\sqrt{\frac{\log |\mathcal{F}| + \log \frac{1}{\delta}}{n}}\right)$
 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 $|\mathcal{F}| = \infty$!