

Empirical Risk Minimization (ERM)

Recap: $(x_1, y_1), \dots, (x_n, y_n)$ - data

$f \in \mathcal{F}$ - candidate hypothesis

$$P_n(f) = \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i)) \quad - \text{empirical risk of } f$$

ERM: choose $\hat{f}_n \in \mathcal{F}$ s.t. $P_n(\hat{f}_n) \leftarrow \min$

$$P_n(\hat{f}_n) = \min_{f \in \mathcal{F}} P_n(f)$$

$$\text{PAC: } \lim_{n \rightarrow \infty} \sup_P P^n \left\{ P(\hat{f}_n) \geq \inf_{f \in \mathcal{F}} P(f) + \varepsilon \right\} = 0 \quad \forall \varepsilon > 0$$

where \hat{f}_n is returned by ERM algo.

Sufficient condition: Uniform Convergence of Empirical Means (UCEM)

$$\lim_{n \rightarrow \infty} \sup_P P^n \left\{ \sup_{f \in \mathcal{F}} |P_n(f) - P(f)| \geq \varepsilon \right\} = 0 \quad \forall \varepsilon > 0$$

(property of l and \mathcal{F} only)

UCEM \Rightarrow ERM is PAC

Notation:

$$(x, y)$$

$$z$$

$$\ell(y, f(x))$$

$$f(z)$$

$$P_n(f) := \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$$

$$P(f) = \mathbb{E}_p[\ell(y, f(x))]$$

$$P_n(f) := \frac{1}{n} \sum_{i=1}^n f(z_i)$$

$$P(f) := \mathbb{E}_p[f(z)]$$

ERM:

$$\hat{f}_n = \underset{f \in \mathcal{F}}{\operatorname{argmin}} P_n(f)$$

$$P_n(A) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{z_i \in A\}}$$

- empirical distribution of z^n

$$z^n = (z_1, \dots, z_n)$$

$$P(A) = P[z \in A]$$

$$\mathbb{E} P_n(A) = P(A) \quad \forall \text{ event } A$$

$$\mathbb{E} P_n(f) = P(f)$$

$$0 \leq f \leq 1 : P\{|P_n(f) - P(f)| \geq \varepsilon\} \leq 2e^{-2n\varepsilon^2}$$

(Hoeffding)

P, P' : two possible prob. dist. for z

$$\begin{aligned} P(f) &:= \mathbb{E}_P[f(z)] \\ P'(f) &:= \mathbb{E}_{P'}[f(z)] \end{aligned} \quad \left. \right\} \quad \forall f \in \mathcal{F}$$

$$\sup_{f \in \mathcal{F}} |P(f) - P'(f)| =: \|P - P'\|_{\mathcal{F}}$$

$\ell^\infty(\mathcal{F})$ seminorm: $\|P - P'\|_{\mathcal{F}}$

$$\leq \|P - P''\|_{\mathcal{F}} + \|P'' - P'\|_{\mathcal{F}}$$

$$\|P - P'\|_{\mathcal{F}} = 0 \Rightarrow P = P'$$

Key quantity : $\|P_n - P\|_{\mathcal{F}}$

$$P_n(\cdot) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{Z_i \in \cdot\}}$$

$$P(\cdot) = P[Z \in \cdot]$$

UCEM: $\|P_n - P\|_{\mathcal{F}} \rightarrow 0$ as $n \rightarrow \infty$

in probability, uniformly in P

Thm (from last lecture - new notation)

Let \hat{f}_n be a random element of \mathcal{F} computed from data Z_1, \dots, Z_n iid P .

$$1) P(\hat{f}_n) \leq P_n(\hat{f}_n) + \|P_n - P\|_{\mathcal{F}}$$

2) if \hat{f}_n minimizes $P_n(\cdot)$ on \mathcal{F} [ERM],
then

$$P(\hat{f}_n) \leq \inf_{f \in \mathcal{F}} P(f) + 2\|P_n - P\|_{\mathcal{F}}.$$

Goal: get a handle on $\|P_n - P\|_{\mathcal{F}}$.

$$\|P_n - P\|_{\mathcal{F}}$$

$$= \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(Z_i) - P(f) \right|$$

$$P_n(f) = \frac{1}{n} \sum_{i=1}^n f(Z_i)$$

Z_1, \dots, Z_n iid P

- f cn of Z_1, \dots, Z_n ,

has bdd differences $c_1 = \dots = c_n = \frac{1}{n}$

$$\Rightarrow P^n \left\{ \|P_n - P\|_{\mathcal{F}} \geq \mathbb{E} \|P_n - P\|_{\mathcal{F}} + \varepsilon \right\} \leq e^{-2n\varepsilon^2}$$

by McDiarmid

$$\text{or : } \|P_n - P\|_{\mathcal{F}} \leq \mathbb{E} \|P_n - P\|_{\mathcal{F}} + \sqrt{\frac{\log(\frac{1}{\delta})}{2n}}$$

w.r.t.
 $\geq 1 - \delta$

How big is $\mathbb{E} \|P_n - P\|_{\mathcal{F}}$?

- symmetrization (Vapnik - Chervonenkis, 1970s;
Giné - Zinn, 1980s)

$$\mathbb{E} P_n(f) = P(f) \quad \forall f$$

$$z_1, \dots, z_n \stackrel{iid}{\sim} P$$

$$\bar{z}_1, \dots, \bar{z}_n \stackrel{iid}{\sim} P, \perp z_1, \dots, z_n \quad (\text{ghost sample})$$

$$\bar{P}_n(f) := \frac{1}{n} \sum_{i=1}^n f(\bar{z}_i)$$

$$P_n(f) \approx P(f), \quad \bar{P}_n(f) \approx P(f) \quad \text{w.h.p.}$$

\rightarrow can expect $P_n(f) - \bar{P}_n(f) \approx 0$ w.h.p.

$$\frac{1}{n} \sum_{i=1}^n \{ f(z_i) - f(\bar{z}_i) \} \quad \begin{matrix} \text{should have} \\ \text{"cancellations"} \end{matrix}$$

$$\|P_n - P\|_{\mathcal{F}} = \sup_{f \in \mathcal{F}} |P_n(f) - \mathbb{E} \bar{P}_n(f)|$$

$$\leq \sup_{f \in \mathcal{F}} \mathbb{E} |P_n(f) - \bar{P}_n(f)|$$

$\mathbb{E} [\cdot] - \text{w.r.t.}$
 $\bar{z}_1, \dots, \bar{z}_n$

$$\leq \mathbb{E} \left\{ \sup_{f \in \mathcal{F}} |P_n(f) - \bar{P}_n(f)| \right\}$$

$$\mathbb{E} \|P_n - P\|_{\mathcal{F}} \leq \mathbb{E} \mathbb{E} \|P_n - \bar{P}_n\|_{\mathcal{F}} \quad z^n, \bar{z}^n \text{ iid}$$

$$\mathbb{E} \|P_n - P\|_{\mathcal{F}} \leq \mathbb{E} \|P_n - \bar{P}_n\|_{\mathcal{F}}$$

$$= \mathbb{E} \left\{ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n [f(z_i) - f(\bar{z}_i)] \right| \right\}$$

$$\forall f_i : f(z_i) - f(\bar{z}_i) \stackrel{d}{=} f(\bar{z}_i) - f(z_i)$$

a r.v. U is symmetric if $U \stackrel{d}{=} -U$

Fact: if U is symmetric, then $U \stackrel{d}{=} \varepsilon U$ where

$\varepsilon \perp U$ and takes values ± 1 w. equal prob.

[ε is a Rademacher r.v.]

\Rightarrow if $\varepsilon_1, \dots, \varepsilon_n \stackrel{iid}{\sim} \text{Rad}$ $\perp\!\!\!\perp z^n, \bar{z}^n$

then $(f(z_i) - f(\bar{z}_i))_i \stackrel{d}{=} (\varepsilon_i (f(z_i) - f(\bar{z}_i)))_{i \in I}$

$$\therefore \mathbb{E} \left\{ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n [f(z_i) - f(\bar{z}_i)] \right| \right\}$$

$$= \mathbb{E} \left\{ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (f(z_i) - f(\bar{z}_i)) \right| \right\} \xrightarrow{\mathbb{E}_{\varepsilon^n, z^n, \bar{z}^n}}$$

$$\leq \mathbb{E} \left\{ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(z_i) \right| \right\} + \mathbb{E} \left\{ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(\bar{z}_i) \right| \right\}$$

$$= 2 \mathbb{E} \left\{ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(z_i) \right| \right\} \xrightarrow{\mathbb{E}_{\varepsilon^n, z^n}}$$

Thm (Gine-Zinn) $\mathbb{E} \|P_n - P\|_{\mathcal{F}} \leq 2 \mathbb{E} \left\{ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(z_i) \right| \right\}$

Rademacher averages (or complexities)

Def A - bdd subset of \mathbb{R}^n

Rademacher average (or complexity) of A :

$$R_n(A) := \mathbb{E} \left\{ \sup_{a \in A} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i a_i \right| \right\}$$

where $a = (a_1, \dots, a_n)^T$ are elements of A , and $\epsilon_1, \dots, \epsilon_n$ are iid Rad r.v.'s

Learning setting:

z_1, \dots, z_n ; \mathcal{F} - fns, $f(z_i) \in [0, 1]$, vi

$$\mathcal{F}(z^n) := \left\{ (f(z_1), \dots, f(z_n)) : f \in \mathcal{F} \right\}$$

- bdd subset of \mathbb{R}^n

$$[\mathcal{F}(z^n) \subseteq [0, 1]^n]$$

$$R_n(\mathcal{F}(z^n)) = \mathbb{E}_{\epsilon} \left\{ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(z_i) \right| \right\}$$

\downarrow
fixed

$$\text{Gine-Zinn: } \mathbb{E} \|P_n - P\|_{\mathcal{F}} \leq 2 \mathbb{E} R_n(\mathcal{F}(z^n))$$