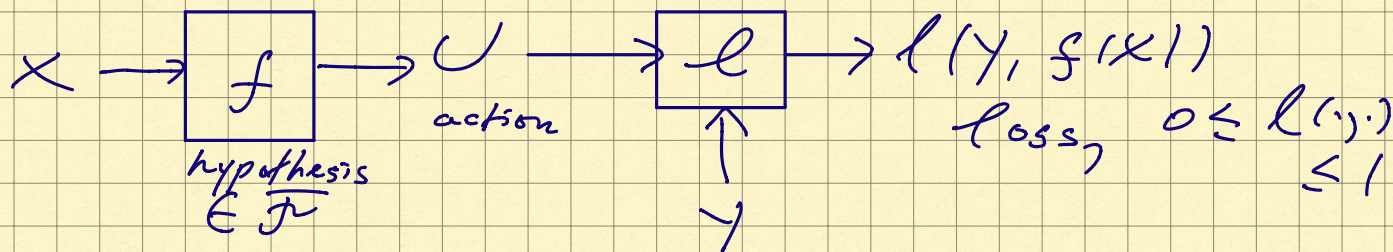


Review: Agnostic (Model-Free) Learning

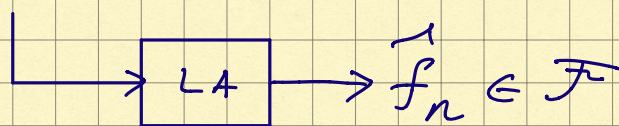
$(X \text{ (feature)}, Y \text{ (label)}) \sim P$ (unmodeled)



$L_P(f) := \mathbb{E}[l(Y, f(X))]$ - risk of f

$L_P^*(\mathcal{F}) := \inf_{f \in \mathcal{F}} L_P(f)$

Learning: $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{iid}{\sim} P$
data



Goal: design LA s.t., no matter what P is,
 $L_P(\hat{f}_n) \approx L_P^*(\mathcal{F})$ w.h.p.

PAC learning: $\forall \epsilon > 0$

$\lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} P^n(L_P(\hat{f}_n) > L_P^*(\mathcal{F}) + \epsilon) = 0$.

Do PAC learning algos exist?

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

Empirical risk of $f \in \mathcal{F}$: $\frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$
 $=: P_n(f)$

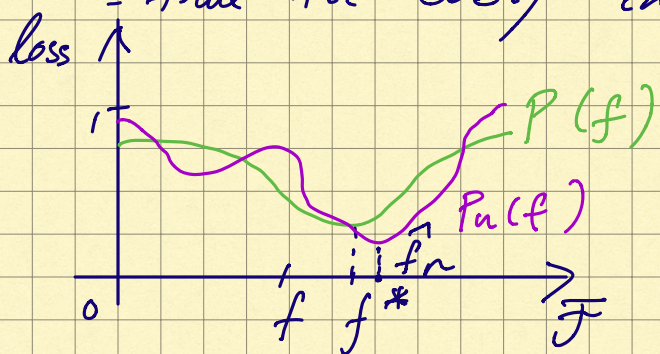
Population risk of $f \in \mathcal{F}$: $L_P(f) = P(f)$

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

2) $|P_n(f) - P(f)| \approx 0$ w.h.p.

$$P\{|P_n(f) - P(f)| > \varepsilon\} \leq 2e^{-2n\varepsilon^2}$$

- true for every individual $f \in \mathcal{F}$



$$P(f^*) = \inf_{f \in \mathcal{F}} P(f)$$

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

Empirical Risk Minimization:

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$$

Mismatched Minimization Lemma

• $G, \hat{G} : \mathcal{U} \rightarrow \mathbb{R}$

• want to minimize G , can't do it directly

• can minimize \hat{G}

• G and \hat{G} are "close": $\sup_{u \in \mathcal{U}} |G(u) - \hat{G}(u)| \leq \varepsilon$

then: 1) $G(u) \leq \widehat{G}(u) + \varepsilon$ for all $u \in \mathcal{U}$
 2) let \hat{u}^* be the minimizer of \widehat{G} ; then

$$G(\hat{u}^*) \leq \inf_{u \in \mathcal{U}} G(u) + 2\varepsilon$$

Proof: 1) is obvious from hypotheses on G, \widehat{G} .
 2) For any $u \in \mathcal{U}$,

$$\begin{aligned} G(u) &\geq \widehat{G}(u) - \varepsilon && |G - \widehat{G}| \leq \varepsilon \\ &\geq \widehat{G}(\hat{u}^*) - \varepsilon && \widehat{G}(u) \geq \widehat{G}(\hat{u}^*) \\ &\geq G(\hat{u}^*) - 2\varepsilon && |G - \widehat{G}| \leq \varepsilon \end{aligned}$$

$\Rightarrow \inf_{u \in \mathcal{U}} G(u) \geq G(\hat{u}^*) - 2\varepsilon.$ □

Our setting: $P, P_n : \mathcal{F} \rightarrow [0, 1]$

if $\sup_{f \in \mathcal{F}} |P_n(f) - P(f)| \leq \varepsilon$, then:

1) for any LA that generates \hat{f}_n ,

$$P(\hat{f}_n) \leq P_n(\hat{f}_n) + \varepsilon$$

2) for ERM, $\hat{f}_n = \operatorname{argmin}_{f \in \mathcal{F}} P_n(f)$,

$$P(\hat{f}_n) \leq \inf_{f \in \mathcal{F}} P(f) + 2\varepsilon.$$

A sufficient condition for ERM to be PAC:

$$\sup_{f \in \mathcal{F}} |P_n(f) - P(f)| \xrightarrow{p} 0 \quad \text{w.h.p.}$$

Formalize: $\mathcal{L}, \mathcal{F}, \mathcal{P}$ given

$$q(n, \varepsilon) := \sup_{P \in \mathcal{P}} P^n \left(\sup_{f \in \mathcal{F}} |P_n(f) - P(f)| \geq \varepsilon \right)$$

ERM is PAC if $q(n, \varepsilon) \rightarrow 0$ as $n \rightarrow \infty$
for any $\varepsilon > 0$.

Example: any finite $\mathcal{F} = \{f_1, \dots, f_M\}$

$$P^n \left(\sup_{f \in \mathcal{F}} |P_n(f) - P(f)| \geq \varepsilon \right)$$

$$= P^n \left(\bigcup_{m=1}^M \{ |P_n(f_m) - P(f_m)| \geq \varepsilon \} \right)$$

$$\leq \sum_{m=1}^M P^n \left(|P_n(f_m) - P(f_m)| \geq \varepsilon \right) \quad (\text{union bound})$$

$$\leq 2M\varepsilon^{-2n} \quad (\text{Hoeffding})$$

$$q(n, \varepsilon) = \sup_P P^n \left(\sup_{f \in \mathcal{F}} |P_n(f) - P(f)| \geq \varepsilon \right) \leq 2M\varepsilon^{-2n} \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Generally: for a broad class of $(\mathcal{L}, \mathcal{F})$,

$$\sup_P \mathbb{E} \left[\sup_{f \in \mathcal{F}} |P_n(f) - P(f)| \right] \leq \frac{C}{\sqrt{n}}$$

where $C > 0$ depends on "complexity" of \mathcal{L}, \mathcal{F}
(\mathcal{F} may be infinite!)

Then ERM IS PAC:

$$P^n \left\{ \sup_{f \in \mathcal{F}} |P_n(f) - P(f)| \geq \varepsilon \right\}$$

$$:= g(z^n)$$

$$z_i = (x_i, y_i)$$

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

Know $\mathbb{E} g(z^n) \leq \frac{C}{\sqrt{n}}$

$$g(z^n) = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)) - \mathbb{E}_P \ell(Y, f(X)) \right|$$

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

$$\Rightarrow P \left\{ g(z^n) \geq \mathbb{E} g(z^n) + t \right\} \leq e^{-2nt^2}$$

$\forall t > 0$ (by McDiarmid)

$$C/\sqrt{n} \leq \varepsilon/2 \quad (\Leftrightarrow) \quad n \geq \frac{4C^2}{\varepsilon^2}$$

$$P \left\{ g(z^n) \geq \varepsilon \right\} = P \left\{ g(z^n) - \mathbb{E} g(z^n) \geq \varepsilon - \mathbb{E} g(z^n) \right\}$$

$$\leq P \left\{ g(z^n) - \mathbb{E} g(z^n) \geq \varepsilon - \frac{C}{\sqrt{n}} \right\}$$

$$\leq P \left\{ g(z^n) - \mathbb{E} g(z^n) \geq \varepsilon/2 \right\}$$

$$\leq e^{-2n(\varepsilon/2)^2}$$

$$= e^{-n\varepsilon^2/2}$$

$$\Rightarrow q(n, \varepsilon) \leq e^{-2n\varepsilon^2/2} \quad \text{for } n \geq \frac{4C^2}{\varepsilon^2}$$