

Review: Concept/Function Learning

(X, Y)

X - features

$Y = f(X)$ - deterministic label
or response

$f: \mathcal{X} \rightarrow \{0, 1\}$ - concept learning

$f: \mathcal{X} \rightarrow [0, 1]$ - function learning

Training data:

$(x_1, f^*(x_1)), (x_2, f^*(x_2)), \dots, (x_n, f^*(x_n))$

where $x_1, \dots, x_n \stackrel{iid}{\sim} P \in \mathcal{P}$

f^* is an element of a given class \mathcal{F}

Learning algo: data $\longrightarrow \hat{f}_n \in \mathcal{F}$

$L_P(f, f^*) := \mathbb{E}_P[|f(X) - f^*(X)|^2]$

$L_P(\hat{f}_n, f^*)$ - random variable, depends on data

Goal: make $L_P(\hat{f}_n) \sim 0$ (with high prob.)

Learning algo is **Probably Approximately Correct (PAC)** if:

$\sup_{P \in \mathcal{P}} \sup_{f \in \mathcal{F}} P_f^n (L_P(\hat{f}_n, f) > \epsilon) \longrightarrow 0$ as $n \rightarrow \infty$

for any $\epsilon > 0$.

Concept learning: \mathcal{C} - family of sets $C \subseteq \mathcal{X}$

$$L_P(C, C^*) = P(C \Delta C^*)$$

$$\mathcal{C} \leftrightarrow \mathcal{F} = \{1_C : C \in \mathcal{C}\}$$

$$1_C(x) = \begin{cases} 1, & \text{if } x \in C \\ 0, & \text{if } x \notin C \end{cases}$$

$$P(C \Delta C^*) = \mathbb{E}_P[|1_C(X) - 1_{C^*}(X)|^2]$$

Limitations:

- 1) deterministic relationship b/w Y and X
 $Y = f(X)$ - what about $Y = f(X) + \text{noise}$
 - 2) the class \mathcal{F} is known - what about mismatch between our assumption on \mathcal{F} and ground truth?
- realizable setting

Agnostic (model-free) formulation (D. Haussler):

- \mathcal{X}, \mathcal{Y} - feature and label space
 $\mathcal{Z} := \mathcal{X} \times \mathcal{Y}$
- \mathcal{U} - action (decision) space
- \mathcal{F} - class of functions $f: \mathcal{X} \rightarrow \mathcal{U}$
[hypothesis space]
- $l: \mathcal{Y} \times \mathcal{U} \rightarrow [0, 1]$ - loss fctn
- \mathcal{P} : class of prob. dist. on $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$

Ideal case: $z = (x, y) \sim P$, known

$$L_P(f) := \mathbb{E}_P[l(y, f(x))]$$

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

Learning! "Nature" selects $P \in \mathcal{P}$

$$\underbrace{(x_1, y_1)}_{z_1}, \dots, \underbrace{(x_n, y_n)}_{z_n} \stackrel{\text{iid}}{\sim} P$$

learning algo $A_n : z^n \rightarrow \hat{f}_n \in \mathcal{F}$

$$L_P(\hat{f}_n) = \int_{\mathcal{X} \times \mathcal{Y}} P(dx, dy) l(y, \hat{f}_n(x))$$

$$0 \leq L_P^*(\mathcal{F}) \leq L_P(\hat{f}_n) \leq 1$$

Goal: construct learning algo s.t.

$$L_P(\hat{f}_n) - L_P^*(\mathcal{F}) \rightarrow 0 \text{ as } n \rightarrow \infty$$

regardless of P .

PAC: for a fixed $\varepsilon > 0$,

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} P^n \left(L_P(\hat{f}_n) \geq L_P^*(\mathcal{F}) + \varepsilon \right) = 0$$

$$:= r_A(n, \varepsilon)$$

A - learning algo.

Ex (function learning)

\mathcal{X} - feature space

$$\mathcal{Y} = \mathcal{U} = [0, 1]$$

\mathcal{F} - target functions

$$\{Y = f(X) \text{ for some } f \in \mathcal{F}\}$$

$$l(y, u) = |y - u|^2$$

\mathcal{P} - all distributions of (X, Y) where $L(X)$ comes from some P_0 (of dist. on \mathcal{X}) and $Y = f(X)$ for some $f \in \mathcal{F}$

$$P_f((X, Y) \in A) = \int_A P(dx) \mathbb{1}_{\{y = f(x)\}}$$

$$\Rightarrow \mathcal{P} = \{P_f : P \in P_0, f \in \mathcal{F}\}$$

$$L_{P_f}(g) = \mathbb{E}_{P_f} |Y - g(X)|^2$$

$$= \int_{\mathcal{X} \times \mathcal{Y}} P(dx) \mathbb{1}_{\{y = f(x)\}} |y - g(x)|^2$$

$$= \int_{\mathcal{X}} P(dx) |f(x) - g(x)|^2 \equiv L_P(f, g)$$

$$L_{P_f}(f) = 0 \quad \Rightarrow \quad L_{P_f}^*(\mathcal{F}) = 0$$

- PAC definitions coincide.

Agnostic case can cover more general situations:

$$Y = f(X) + W$$

$$\mathbb{E}W = 0, \text{ indep. of } X$$

f may not be in \mathcal{F}

When is PAC learning possible in model-free case?

$$(x_1, y_1), \dots, (x_n, y_n) \stackrel{\text{iid}}{\sim} P$$

$$l_f : \mathcal{Z} \rightarrow [0, 1] \quad : \quad l_f(z) = l_f(x, y) \\ = l(y, f(x))$$

$$P(l_f) := \int P(dx, dy) l_f(x, y) \quad - \text{expected loss of } f$$

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

- can be computed for every $f \in \mathcal{F}$

$$P \left\{ \left| P_n(l_f) - P(l_f) \right| \geq \epsilon \right\} \leq 2e^{-2n\epsilon^2} \\ \text{(Hoeffding)}$$

- in particular, if $f^* \in \mathcal{F}$ achieves $L_P^*(\mathcal{F})$, then

$$P \left\{ \left| P_n(l_{f^*}) - L_P^*(\mathcal{F}) \right| \geq \epsilon \right\} \leq 2e^{-2n\epsilon^2}$$

Candidate learning algo: Empirical Risk Minimization (ERM)

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

Claim: (to be proved in next lecture)

if l (loss), P , and \mathcal{F} (hypothesis class) are such that

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} P^n \left(\sup_{f \in \mathcal{F}} |P_n(l_f) - P(l_f)| \geq \epsilon \right) = 0$$

for any $\epsilon > 0$, then ERM is a PAC learning algo.

Exercise: prove that this holds if $|\mathcal{F}| < \infty$.