

Generalization Theory

2020年10月12日 14:15

→ there is no universal learner.

THEOREM 5.1 (No-Free-Lunch) Let A be any learning algorithm for the task of binary classification with respect to the 0-1 loss over a domain \mathcal{X} . Let m be any number smaller than $|\mathcal{X}|/2$, representing a training set size. Then, there exists a distribution \mathcal{D} over $\mathcal{X} \times \{0, 1\}$ such that:

Here the hypothesis class is the whole function space!

1. There exists a function $f : \mathcal{X} \rightarrow \{0, 1\}$ with $L_{\mathcal{D}}(f) = 0$.
2. With probability of at least $1/7$ over the choice of $S \sim \mathcal{D}^m$ we have that $L_{\mathcal{D}}(A(S)) \geq 1/8$.

Proof. Assume $m = |\mathcal{X}|/2$
 $T = 2^{2^m}$ functions $\mathcal{X} \rightarrow \{0, 1\}$

All possible functions f — distribution \mathcal{D}
 $\Rightarrow L_{\mathcal{D}}(f) = 0$

$\{f_1, \dots, f_T\}$.

def. $D_i(\{x, y\}) = \begin{cases} \frac{1}{|\mathcal{X}|} & \text{if } y = f_i(x) \\ 0 & \text{otherwise} \end{cases}$ (distribution)

All kinds of sampled training set S
 bound the expected loss $\max_i \mathbb{E}_{S \sim D_i} [L_{D_i}(A(S))]$ by $\frac{1}{4}$
 by avg loss of f class

$\Rightarrow L_{D_i}(f_i) = 0$.

There are $k = (2^m)^m$ kinds of sequence of sampling result $\Rightarrow \{S_1, \dots, S_k\}$

$S_j = (x_1, \dots, x_m)$

$S_j^i = (x_1, f_i(x_1), \dots, (x_m, f_i(x_m)))$

$\Rightarrow \mathbb{E}_{S \sim D_i} [L_{D_i}(A(S))] = \frac{1}{k} \sum_{j=1}^k L_{D_i}(A(S_j^i))$
 function derived by algo. A .

$\max_{i \in [T]} \frac{1}{k} \sum_{j=1}^k L_{D_i}(A(S_j^i)) \geq \frac{1}{T} \sum_{i=1}^T \frac{1}{k} \sum_{j=1}^k L_{D_i}(A(S_j^i))$
 $\max \geq \text{Avg} \geq \min \geq \frac{1}{T} \sum_{i=1}^T \frac{1}{k} \sum_{j=1}^k L_{D_i}(A(S_j^i))$

Let v_1, \dots, v_p be examples in $C = S_j$

$\Rightarrow p \geq m$

$\Rightarrow \forall$ function h

$$\begin{aligned} L_{D_i}(h) &= \frac{1}{2^m} \sum_{x \in C} [h(x) \neq f_i(x)] \\ &\geq \frac{1}{2^m} \sum_{r=1}^p [h(v_r) \neq f_i(v_r)] \\ &\geq \frac{1}{2^p} \sum_{r=1}^p [h(v_r) \neq f_i(v_r)] \end{aligned}$$

$$\begin{aligned} \Rightarrow \frac{1}{T} \sum_{i=1}^T L_{D_i}(A(S_j^i)) &\geq \frac{1}{T} \sum_{i=1}^T \frac{1}{2^p} \sum_{r=1}^p [A(S_j^i)(v_r) \neq f_i(v_r)] \\ &= \frac{1}{2^p} \sum_{r=1}^p \frac{1}{T} \sum_{i=1}^T [A(S_j^i)(v_r) \neq f_i(v_r)] \end{aligned}$$

$$\begin{aligned} \Rightarrow \frac{1}{2} \sum_{i=1}^p \sum_{j=1}^p [A(S_j^i)(V_i) \neq f_i(V_i)] \\ = \frac{1}{2p} \sum_{i=1}^p \frac{1}{p} \sum_{j=1}^p [A(S_j^i)(V_i) \neq f_i(V_i)] \\ \geq \frac{1}{2} \cdot \min_{\text{recp}} \frac{1}{p} \sum_{i=1}^p [A(S_j^i)(V_i) \neq f_i(V_i)] \end{aligned}$$

partition T functions into $\frac{1}{2}$ pairs. (f_i, f_i') only different on $f_i(V_i), f_i'(V_i)$.

$$\Rightarrow \frac{1}{p} \sum_{i=1}^p [A(S_j^i)(V_i) \neq f_i(V_i)] = \frac{1}{2}$$

$$\Rightarrow \left[\begin{array}{l} \forall \text{ algo } A', \exists f: X \rightarrow \{0,1\} \text{ and distribution } D. \\ \text{st. } \mathbb{E}_{S \sim D^n} [L_D(A'(S))] \geq \frac{1}{4} \end{array} \right]$$

Hypothesis Class

- the class of functions

- infinite / finite

eg. \downarrow \swarrow eg. decision tree

$$H = \{f: f(x) = w^T x, \|w\|_2 \leq 1\}$$

Empirical risk minimization.

$$\hookrightarrow \text{ERM}_H(S) \in \underset{h \in H}{\text{argmin}} L_S(h)$$

\uparrow
training dataset

If H is infinite, ERM_H may simply memorize S (overfit)

Def. 2.1 Realizability Assumption

There exists $h^* \in H$, st. $L_{(D,f)}(h^*) = 0$ h : hypo. function.
 $\forall S$ sampled from D , labeled by f , $L_S(h^*) = 0$

Corollary 2.3 If finite hypothesis class, $\delta \in (0,1), \epsilon > 0$

$$\text{let integer } m: m \geq \frac{\log(|H|/\delta)}{\epsilon}$$

Then $\forall D, f$ that realizability assumption holds.

w.h.p of $1-\delta$, over i.i.d. sample S of size m , we have

$$\forall \text{ERM hypo. } h_S: L_{(D,f)}(h_S) < \epsilon$$

Proof. want to upper bound:

$$\bar{S} = \{S: L_{(D,f)}(h_S) > \epsilon\}$$

$$\text{bad hypothesis: } H_B = \{h \in H \mid L_{(D,f)}(h) > \epsilon\}$$

$$\text{misleading samples: } \mathcal{M} = \{S: \exists h \in H_B, L_S(h) = 0\}$$

bad hypothesis: $\mathcal{H}_B = \{h \in \mathcal{H} \mid L_{D,f}(h) > \epsilon\}$

misleading samples: $\mathcal{M} = \{S: \exists h \in \mathcal{H}_B. L_S(h) = 0\}$

By realization assumption, $L_S(h) = 0$

$\Rightarrow L_{D,f}(h) > \epsilon$ only if S is misleading

$$\Rightarrow \bar{\mathcal{S}} \subseteq \mathcal{M} = \bigcup_{h \in \mathcal{H}_B} \{S \mid L_S(h) = 0\}$$

$$D^m(\bar{\mathcal{S}}) \leq \sum_{h \in \mathcal{H}_B} D^m(\{S \mid L_S(h) = 0\}) \quad \text{Union Bound}$$

$\downarrow S = \{x_1, \dots, x_m\}$

$$= \sum_{h \in \mathcal{H}_B} \prod_{i=1}^m D(\{x_i: h(x_i) = f(x_i)\})$$

$$= \sum_{h \in \mathcal{H}_B} \prod_{i=1}^m (1 - L_{D,f}(h))$$

$$\leq \sum_{h \in \mathcal{H}_B} \prod_{i=1}^m (1 - \epsilon)$$

$$\leq \sum_{h \in \mathcal{H}_B} e^{-\epsilon m} \leq |\mathcal{H}| e^{-\epsilon m}$$

Def 3.1 PAC Learnability

A hypo. class \mathcal{H} is PAC learnable if:

\exists function $m_{\mathcal{H}}: (0, 1]^2 \rightarrow \mathbb{N}$, Learning algo. A .

$\forall D, f$ if realization assumption holds, run A on $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ examples (iid)

$\Rightarrow A$ returns h w.p. $\geq 1 - \delta$. $L_{D,f}(h) \leq \epsilon$

e.g. ERM for finite \mathcal{H} is PAC-learnable

Realizability too strong - Sometimes $L_{D,f}(h^*)$ cannot reach 0

- f may have multiple labels.

(weaker version)

the Bayes optimal predictor

$$f_D(x) = \begin{cases} 1 & \text{if } \Pr[y=1|x] \geq \frac{1}{2} \\ 0 & \text{otherwise.} \end{cases}$$

\Rightarrow expect $L_D(h) \leq \min_{h \in \mathcal{H}} L_D(h) + \epsilon$ (instead of $\leq \epsilon$)

when realizability holds. \rightarrow
degrade to

\downarrow use this to replace the requirement in PAC-learnability

Agnostic PAC Learnability

Error Decomposition

$$L_D(h_S) = \epsilon_{\text{app}} + \epsilon_{\text{est}}$$

- ...

$$L_D(h_S) = \epsilon_{app} + \epsilon_{est}$$

$$\left\{ \begin{array}{l} \epsilon_{app}: \text{approximation} = \min_{h \in \mathcal{H}} L_D(h) = \underbrace{L_D(BD)}_{\text{base opt. due to inheritance of } D} + \epsilon'_{app} \\ \epsilon_{est}: \text{estimation} \end{array} \right.$$

Infinite \mathcal{H} ?

Lema 6.2. \mathcal{H} : class of thresholds (1D)

$$\text{ERM rule \& } m_{\mathcal{H}}(\epsilon, \delta) \leq \lceil \log(2/\delta)/\epsilon \rceil \quad \text{PAC learnable.}$$

Generally \Rightarrow VC dimension

Def. Restriction of \mathcal{H} to C

Let $C = \{c_1, \dots, c_m\} \subset X$

$$\Rightarrow \mathcal{H}_C = \{h(c_1), \dots, h(c_m) : h \in \mathcal{H}\}$$

Def. If all function $h: C \rightarrow \{0, 1\}$ is in \mathcal{H}_C , then \mathcal{H} shatters C .

Def. VC dimension

$\text{VCdim}(\mathcal{H})$: maximal size of $C \subset X$ that can be shattered by \mathcal{H} .

THEOREM 6.8 (The Fundamental Theorem of Statistical Learning – Quantitative Version) Let \mathcal{H} be a hypothesis class of functions from a domain X to $\{0, 1\}$ and let the loss function be the 0 – 1 loss. Assume that $\text{VCdim}(\mathcal{H}) = d < \infty$. Then, there are absolute constants C_1, C_2 such that:

1. \mathcal{H} has the uniform convergence property with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\epsilon^2} \leq m_{\mathcal{H}}^{uc}(\epsilon, \delta) \leq C_2 \frac{d + \log(1/\delta)}{\epsilon^2}$$

2. \mathcal{H} is agnostic PAC learnable with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\epsilon^2} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq C_2 \frac{d + \log(1/\delta)}{\epsilon^2}$$

3. \mathcal{H} is PAC learnable with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\epsilon} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq C_2 \frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon}$$

Def. ϵ -Representative Sample

A training set S .

$$\sup_{h \in \mathcal{H}} |L_D(h) - L_S(h)| \leq \epsilon.$$

$\mathcal{F} = \ell \circ \mathcal{H}$ (loss, hypo)

Def. Representativeness of S w.r.t. \mathcal{F} :

$$\text{Rep.}(\mathcal{F}, S) \stackrel{\text{def}}{=} \sup (L_D(f) - L_S(f))$$

Def. Representativeness of S w.r.t. F :

$$\text{Rep}_D(F, S) \stackrel{\text{def}}{=} \sup_{f \in F} (L_D(f) - L_S(f))$$

\uparrow
 $= \mathbb{E} f(z)$
 $z \sim D$

Split S into two parts to estimate this:

Rademacher variable $\sigma_i = \begin{cases} \frac{1}{2} & p=0.5 \\ -\frac{1}{2} & p=0.5 \end{cases}$

Rademacher complexity of F w.r.t. S :

$$R(F \circ S) \stackrel{\text{def}}{=} \frac{1}{m} \mathbb{E}_{\sigma \sim \{\pm 1\}^m} \left[\sup_{f \in F} \sum_{i=1}^m \sigma_i f(z_i) \right]$$

Lemma.

$$\mathbb{E}_{S \sim D^m} [\text{Rep}_D(F, S)] \leq 2 \mathbb{E}_{S \sim D^m} R(F \circ S)$$

Proof. $L_D(f) - L_S(f) = \mathbb{E}_{S'} (L_{S'}(f) - L_S(f))$

$$\begin{aligned} \mathbb{E}_S \left[\sup_f (L_D(f) - L_S(f)) \right] &\leq \mathbb{E}_S \left[\mathbb{E}_{S'} \left[\sup_f (L_{S'}(f) - L_S(f)) \right] \right] \\ &= \frac{1}{m} \mathbb{E} \left[\sup_{f \in F} (f(z'_1) - f(z_1)) \right] \end{aligned}$$

z'_j, z_j are iid & indep.

$$\begin{aligned} \mathbb{E}_{S, S', f} \sup_f \left[f(z'_j) - f(z_j) + \sum_{i \neq j} f(z'_i) - f(z_i) \right] &\quad \textcircled{1} \\ = \mathbb{E}_{S, S', f} \sup_f \left[f(z_j) - f(z'_j) + \sum_{i \neq j} f(z'_i) - f(z_i) \right] &\quad \textcircled{2} \end{aligned}$$

$$\begin{aligned} &\Rightarrow \mathbb{E}_{S, S', \sigma_j} \sup_f \left[(f(z'_j) - f(z_j)) \sigma_j + \sum_{i \neq j} f(z'_i) - f(z_i) \right] \\ &= \frac{1}{2} \textcircled{1} + \frac{1}{2} \textcircled{2} \\ &= \textcircled{1} \end{aligned}$$

$$\begin{aligned} \hookrightarrow \mathbb{E}_{S, S', f} \sup_f \left(\sum_{i=1}^m f(z'_i) - f(z_i) \right) &= \mathbb{E}_{S, S', \sigma} \sup_f \left[\sum_{i=1}^m \sigma_i (f(z'_i) - f(z_i)) \right] \\ &\downarrow \\ &\sup_f \left[\sum_{i=1}^m \sigma_i (f(z'_i) - f(z_i)) \right] \\ &\leq \sup_f \sum_{i=1}^m \sigma_i f(z'_i) + \sup_f \sum_{i=1}^m \sigma_i f(z_i) \end{aligned}$$

$$\begin{aligned} \Rightarrow \mathbb{E}_S [\text{Rep}_D(F, S)] &\leq \frac{2}{m} \mathbb{E}_{S, S', \sigma} \left[\sup_f \sum_{i=1}^m \sigma_i f(z_i) \right] \\ &= 2 \mathbb{E}_S [R(F \circ S)] \end{aligned}$$

THEOREM 26.5 Assume that for all z and $h \in \mathcal{H}$ we have that $|\ell(h, z)| \leq c$.
Then,

1. With probability of at least $1 - \delta$, for all $h \in \mathcal{H}$,

$$L_{\mathcal{D}}(h) - L_S(h) \leq 2 \mathbb{E}_{S' \sim D^m} R(\ell \circ \mathcal{H} \circ S') + c \sqrt{\frac{2 \ln(2/\delta)}{m}}.$$

In particular, this holds for $h = \text{ERM}_{\mathcal{H}}(S)$.

2. With probability of at least $1 - \delta$, for all $h \in \mathcal{H}$,

$$L_{\mathcal{D}}(h) - L_S(h) \leq 2 R(\ell \circ \mathcal{H} \circ S) + 4c \sqrt{\frac{2 \ln(4/\delta)}{m}}.$$

In particular, this holds for $h = \text{ERM}_{\mathcal{H}}(S)$.

3. For any h^* , with probability of at least $1 - \delta$,

$$L_{\mathcal{D}}(\text{ERM}_{\mathcal{H}}(S)) - L_{\mathcal{D}}(h^*) \leq 2 R(\ell \circ \mathcal{H} \circ S) + 5c \sqrt{\frac{2 \ln(8/\delta)}{m}}.$$