

# Review of Chapter07

## Understanding Machine Learning

Peng Lingwei

July 17, 2019

# Uniform Convergence

## Definition

A hypothesis class  $\mathcal{H}$  is *Uniform Convergence* if:

$\exists m_{\mathcal{H}}^{UC}(\delta, \epsilon) \rightarrow \mathbb{N}$  satisfies:

For  $\forall \epsilon, \delta \in (0, 1)$ ,

The training set  $\{S : |S| \geq m_{\mathcal{H}}^{UC}(\epsilon, \delta), S \sim \mathcal{D}^m\}$  guarantees that  $\mathbb{P}\{\forall h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon\} \geq 1 - \delta$ .

## Theorem

$$m_{\mathcal{H}}^{UC}(\epsilon, \delta) \leq \left\lceil \frac{\log(2|\mathcal{H}|/\delta)}{2\epsilon^2} \right\rceil$$

## Proof.

1.  $L_{S=\{(x_1, y_1), \dots, (x_m, y_m)\} \sim \mathcal{D}^m}(h) = \frac{1}{m} \sum_{i=1}^m 1\{h(x_i) \neq y_i\}$
2.  $L_{\mathcal{D}} = \mathbb{E}_{(x, y) \sim \mathcal{D}}\{1\{h(x) \neq y\}\} = \mathbb{E}_{S \sim \mathcal{D}^m}\{L_S(h)\}$
3.  $\mathbb{P}\{\exists h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| \geq \epsilon\} \leq \sum_{i=1}^{|\mathcal{H}|} \mathbb{P}\{h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| \geq \epsilon\} \leq 2|\mathcal{H}| \exp(-2m\epsilon^2)$

# Nonuniformly Learnable

## Definition

A hypothesis class  $\mathcal{H}$  is *Nonuniform Learnable* if:

$\exists A : S \rightarrow h_S \in \mathcal{H}, m_{\mathcal{H}}^{NUL}(\delta, \epsilon, h) \rightarrow \mathbb{N}$  satisfies:

For  $\forall \epsilon, \delta \in (0, 1), h \in \mathcal{H}$ ,

The training set  $\{S : |S| \geq m_{\mathcal{H}}^{NUL}(\epsilon, \delta, h), S \sim \mathcal{D}^m\}$  guarantees that

$$\mathbb{P}\{L_{\mathcal{D}}(A(S)) - L_{\mathcal{D}}(h) \leq \epsilon\} \geq 1 - \delta$$

## Theorem

$\mathcal{H} = \cup_{n \in \mathbb{N}} \mathcal{H}_n$  (means countable sets' union), s.t.  $\mathcal{H}_n$  is uniform convergence.

$\Rightarrow \mathcal{H}$  is nonuniformly learnable.

## Theorem

$\mathcal{H}$  of binary classifiers is nonuniformly learnable.

$\iff \mathcal{H} = \cup_{n \in \mathbb{N}} \mathcal{H}_n$ , s.t.  $\mathcal{H}_n$  is agnostic PAC learnable.

# Structural Risk Minimization

## Definition

$$\epsilon_n(m, \delta) = \min\{\epsilon \in (0, 1) : m_{\mathcal{H}_n}^{UC}(\epsilon, \delta) \leq m\}$$

## Theorem

$$\mathbb{P}\{\forall h \in \mathcal{H}, L_{\mathcal{D}}(h) - L_{\mathcal{S}}(h) \leq \min_{n: h \in \mathcal{H}_n} \epsilon_n(m, w(n)\delta)\} \geq 1 - \delta$$

## Proof.

$$\mathbb{P}\{\forall h \in \mathcal{H}_n, |L_{\mathcal{S}}(h) - L_{\mathcal{D}}(h)| \leq \epsilon_n(m, w(n)\delta)\} \geq 1 - w(n)\delta$$

$$\mathbb{P}\{\exists h \in \mathcal{H}_n, |L_{\mathcal{S}}(h) - L_{\mathcal{D}}(h)| \geq \epsilon_n(m, w(n)\delta)\} \leq w(n)\delta$$

$$\mathbb{P}\{\exists h \in \mathcal{H}, |L_{\mathcal{S}}(h) - L_{\mathcal{D}}(h)| \geq \epsilon_{n: h \in \mathcal{H}_n}(m, w(n)\delta)\} \leq \sum_n w(n)\delta \leq \delta$$

$$\mathbb{P}\{\forall h \in \mathcal{H}, |L_{\mathcal{S}}(h) - L_{\mathcal{D}}(h)| \leq \epsilon_{n: h \in \mathcal{H}_n}(m, w(n)\delta)\} \leq 1 - \delta$$

$$\mathbb{P}\{\forall h \in \mathcal{H}, L_{\mathcal{D}}(h) - L_{\mathcal{S}}(h) \leq \min_{n: h \in \mathcal{H}_n} \epsilon_n(m, w(n)\delta)\} \leq 1 - \delta \quad \square$$

## Definition

(Structural Risk Minimization)

1. **prior knowledge:**

- ▶  $\mathcal{H} = \cup_n \mathcal{H}_n$  where  $\mathcal{H}_n$  has uniform convergence with  $m_{\mathcal{H}_n}^{UC}$ .
- ▶  $w : \mathbb{N} \rightarrow [0, 1]$  s.t.  $\sum_n w(n) \leq 1$

2. **define:**  $\epsilon_n$  and  $n(h) = \min\{n : h \in \mathcal{H}_n\}$

3. **input:** training set  $S \sim \mathcal{D}^m$ , confidence  $\delta$

4. **output:**  $h \in \arg \min_{h \in \mathcal{H}} [L_S(h) + \epsilon_{n(h)}(m, w(n(h)))\delta]$

## Theorem

$$m_{\mathcal{H}}^{NUL}(\epsilon, \delta, h) \leq m_{\mathcal{H}_{n(h)}}^{UC} \left( \epsilon/2, \frac{6\delta/2}{(\pi n(h))^2} \right)$$

### Proof.

Let  $m \geq m_{\mathcal{H}_{n(h)}}^{UC} \left( \epsilon/2, \frac{6\delta/2}{(\pi n(h))^2} \right)$ , then

$$\mathbb{P}\{\forall h \in \mathcal{H}, L_{\mathcal{D}}(h) - L_S(h) \leq \epsilon_n(m, w(n(h))\delta)\} \geq 1 - \delta/2$$

$$\mathbb{P}\{\forall h \in \mathcal{H}, L_{\mathcal{D}}(A(S)) \leq L_S(h) + \epsilon/2\} \geq 1 - \delta/2.$$

From uniform convergence property, we also can get:

$\mathbb{P}\{\forall h \in \mathcal{H}, L_S(h) \leq L_D(h) + \epsilon/2\} \geq 1 - \delta/2$  Then we can guarantee nonuniformly learnable event happens:

$$\mathbb{P}\{\forall h \in \mathcal{H}, L_{\mathcal{D}}(A(S)) \leq L_D(h) + \epsilon\} \geq 1 - \delta$$



In chapter6:

If  $VCdim(\mathcal{H}) = n$ , then  $m_{\mathcal{H}_n}^{UC}(\epsilon, \delta) = C \frac{n + \log(1/\delta)}{\epsilon^2}$ .

Theorem

$$\begin{aligned} m_{\mathcal{H}}^{NUL}(\epsilon, \delta, h) - m_{\mathcal{H}_n}^{UC}(\epsilon/2, \delta) &\leq m_{\mathcal{H}_n}^{UC}\left(\epsilon/2, \frac{3\delta}{(\pi n)^2}\right) - m_{\mathcal{H}_n}^{UC}(\epsilon/2, \delta) \\ &\leq \frac{4C}{\epsilon^2} \log\left(\frac{(\pi n)^2}{3}\right) \\ &\leq \frac{4C}{\epsilon^2} 2 \log\left(\frac{\pi n}{\sqrt{3}}\right) \\ &\leq \frac{4C}{\epsilon^2} 2 \log(2n) \end{aligned}$$

# Consistency

## Definition

A hypothesis class  $\mathcal{H}$  is *Consistency* in probability distributions set  $\mathcal{P}$  if:

$\exists A : S \rightarrow h_S \in \mathcal{H}, m_{\mathcal{H}}^{CON}(\delta, \epsilon, h, \mathcal{D}) \rightarrow \mathbb{N}$  satisfies:

For  $\forall \epsilon, \delta \in (0, 1), h \in \mathcal{H}, \mathcal{D} \in \mathcal{P}$ ,

The training set  $\{S : |S| \geq m_{\mathcal{H}}^{CON}(\epsilon, \delta, h), S \sim \mathcal{D}^m\}$  guarantees that

$$\mathbb{P}\{L_{\mathcal{D}}(A(S)) - L_{\mathcal{D}}(h) \leq \epsilon\} \geq 1 - \delta$$



## Theorem

The algorithm **Memorize** is consistency in countable  $\mathcal{X}$ . ( $\mathcal{P}$  is the set of every distribution on  $\mathcal{X}$ )

## Proof.

Let an  $\mathcal{X}$ 's enumeration  $\{x_i : i \in \mathbb{N}\}$  satisfies:

$$i \leq j \Leftrightarrow \mathcal{D}(x_i) \geq \mathcal{D}(x_j).$$

It's easy to verify  $\lim_{n \rightarrow \infty} \sum_{i=n}^{\infty} \mathcal{D}(x_i) = 0$ , which means that

$$\forall \epsilon > 0, \exists N \in \mathbb{N}, \text{ such that } \mathcal{D}(i > N, x_i) < \epsilon$$

$$\begin{aligned} \mathbb{P}\{\exists x \notin S, \mathcal{D}(x) \geq \epsilon\} &\leq \mathbb{P}\{\exists i \leq N, x_i \notin S\} \\ &\leq \sum_{i=1}^N \mathbb{P}\{x_i \notin S\} = \sum_{i=1}^N (1 - \mathcal{D}(x_i))^m \\ &\leq N(1 - \epsilon)^m \leq Ne^{-\epsilon m} \end{aligned}$$