

Non-uniform learning

- If $VC(H)$ is finite, \widehat{err}_S converges uniformly to err_D .

That is, with probability at least $1-\delta$,

$$\sup_{h \in H} |\widehat{err}_S(h) - err_D(h)| \leq \epsilon$$

assuming $|S| \geq \tilde{\Omega}\left(\frac{VC(H) + \log(1/\delta)}{\epsilon^2}\right)$

- Note that the guarantee

$$|\widehat{err}_S(h) - err_D(h)|$$

does NOT depend on n .

- That is, the convergence is uniform over all $h \in H$.

- We now consider non-uniform convergence.

Definition: A hypothesis class

$H \subseteq Y^X$ is non-uniformly learnable

iff there exists a learning algorithm A such that for all $\epsilon, \delta \in (0, 1)$, $h \in H$ there exists positive integer m such that for every distribution D over $X \times Y$ if S is an i.i.d. sample from D of size m

then

$$\text{err}_D(A(S)) \leq \text{err}_D(h) + \epsilon.$$

1) Note the order of quantifiers:

$$\bullet \forall \epsilon \forall \delta \forall h \in H \exists m \forall D \dots$$

Compare that with agnostic PAC:

$$\bullet \forall \epsilon \forall \delta \exists m \forall D \forall h \dots$$

2) Note that if H is agnostically PAC learnable then it is non-uniformly learnable.

3) The reverse implication is false.

4) Given a learning algorithm A , we define $m_A(\epsilon, \delta, h)$ to be the smallest sample size such that if $|S| \geq m_A(\epsilon, \delta, h)$ then with probability at least $1 - \delta$,

$$\text{err}_D(A(S)) \leq \text{err}_D(h) + \epsilon.$$

Lemma: If H is uniformly learnable then H is a countable union of classes with finite VC-dimension. That is,

$$H = \bigcup_{h=1}^{\infty} H_h$$

... ..

where $VC(H_n) < +\infty$ for all n .

Proof: • Let A be a learning algorithm for non-uniform learning of H .

• Let $m_A : (0, 1)^2 \times H \rightarrow \mathbb{N}$ be the sample size function.

• Define

$$H_n = \{ h \in H : m_A(1/8, 1/7, h) \leq n \}$$

• Clearly

$$H = \bigcup_{n=1}^{\infty} H_n$$

... .. 1

• H_n is PAC learnable using A
since $m_A(1/8, 1/7, n) \leq n$.

• Thus H_n has finite VC dimension.

Example:

$$H_n = \{ \text{sign}(p(x)) : p(x) = a_n x^n + \dots + a_1 x + a_0 \}$$

$$H = \{ \text{sign}(p(x)) : p(x) \text{ is a polynomial} \}$$

• H_n is (agnostically) PAC learnable

• H is NOT (agnostically) PAC learn.

-
- We will see later that H is non-uniformly learnable.
-

Structural risk minimization

- Suppose $H = \bigcup_{n=1}^{\infty} H_n$ where

H_1, H_2, \dots have finite VC-dims.

- Each H_n satisfies w.p. $1-\delta$,

$$\sup_{h \in H_n} |\widehat{\text{err}}_S(h) - \text{err}_D(h)| \leq \epsilon$$

if $|S| \geq m_{\text{lin}}(\epsilon, \delta, VC(H_n))$

• Consider $w: \mathbb{N} \rightarrow [0,1]$ such that

$$\sum_{n=1}^{\infty} w(n) = 1$$

• w can be thought of as

a) weight function

b) probability distribution

c) prior distribution

• For example

$$w(n) = 2^{-n} \quad \text{or} \quad w(n) = \frac{6}{\pi^2 n^2}$$

• If H is finite union $H_1 \cup \dots \cup H_N$,

$$w(n) = \frac{1}{N}$$

• Let

$$\varepsilon_n(m, \delta) = \min \left\{ \varepsilon \in (0, 1) : m_{\text{ve}}(\varepsilon, \delta, \text{VC}(H_n)) \leq m \right\}$$

Theorem:

Let D be any distribution over $X \times Y$. Let $H = \bigcup_{n=1}^{\infty} H_n$ where

$\text{VC}(H_n) < +\infty$. Let $w: \mathbb{N} \rightarrow [0, 1]$

s.t. $\sum_{n=1}^{\infty} w(n) = 1$. Let m be

a positive integer. Let $\delta \in (0, 1)$.

Let S be i.i.d. sample from

\mathbb{N} of size m . Then, with

probability at least $1-\delta$,
for all $n \geq 1$ and all $h \in H_n$,

$$|\widehat{\text{err}}_S(h) - \text{err}_D(h)| \leq \epsilon_n(m, \omega(n) \cdot \delta).$$

Also, w.p. $1-\delta$, for all $h \in H$,

$$\text{err}_D(h) \leq \widehat{\text{err}}_S(h) + \min_{n: h \in H_n} \epsilon_n(m, \omega(n) \cdot \delta)$$

Proof:

• W.p. $1-\delta \cdot \omega(n)$

$$\sup_{h \in H_n} |\text{err}_D(h) - \widehat{\text{err}}_S(h)| \leq \epsilon_n(m, \omega(n) \cdot \delta)$$

|| || || || ||

- Union bound over all $m = 1, 2, \dots$ implies that w.p. $1 - \delta$ for all $n \geq 1$,

$$\sup_{h \in H_n} |\text{err}_D(h) - \widehat{\text{err}}_S(h)| \leq \varepsilon_n(m, w(n) \cdot \delta)$$

- The second part is trivial consequence of the first one.
-

- Let $n(h) = \min \{n \in \mathcal{N} : h \in H_n\}$

- Theorem implies that

$$\text{err}(h) \leq \widehat{\text{err}}_S(h) + \varepsilon_{n(h)}(m, w(n(h)) \cdot \delta)$$

$$\text{err}_D(h) = \underbrace{\dots}_{\text{Can be computed from } S \text{ only !!!}}$$

Can be computed from
 S only !!!

SRM algorithm

Parameters: $\delta \in (0, 1)$, decomposition
 $H = \bigcup_{n=1}^{\infty} H_n$, $w: \mathbb{N} \rightarrow [0, 1]$ s.t. $\sum_{n=1}^{\infty} w(n) = 1$

Input: Labeled sample S

Output:

$$\hat{h} = \underset{h \in H}{\text{argmin}} \widehat{\text{err}}_S(h) + E_n(h) (m, w(n) \cdot \delta)$$

Lemma:

Let H_1, H_2, \dots be classes with finite VC-dimension. Let

$H = \bigcup_{n=1}^{\infty} H_n$. Then H is

non-uniformly learnable using

SRTM with, say, $w(n) = 2^{-n}$

and, say, $\delta = 1/|S|$.

Proof:

• Let $\epsilon, \delta \in (0, 1)$ and $h \in H$.

• Consider smallest n such that

$$h \in H_n.$$

• Let $m \geq \max \left\{ \frac{1}{\delta}, m_{VC} \left(\epsilon/2, \frac{\delta}{2^n}, VC(H_n) \right) \right\}$

• Let S be i.i.d. sample of size m .

R . $\exists h' \in H$

- By previous

$$|\text{err}_D(h') - \widehat{\text{err}}_S(h')| \leq \mathcal{E}_{u(h)}(m, \omega(u(h'))) \cdot \delta$$

w.p. at least $1 - \delta$

- Let $\hat{h} = \text{SRM}(S)$

- By definition of SRM

$$\widehat{\text{err}}_S(\hat{h}) + \mathcal{E}_{u(h)}(m, \omega(u(\hat{h}))) \cdot \frac{1}{m} \leq$$

$$\widehat{\text{err}}_S(h) + \mathcal{E}_{u(h)}(m, \omega(u(h))) \cdot \frac{1}{m}.$$

- Therefore, w.p. $1 - \delta$,

$$\text{err}_D(\hat{h}) \leq \widehat{\text{err}}_S(h) + \mathcal{E}_{u(h)}(m, \omega(u(h))) \cdot \frac{1}{m}$$

$$\leq \widehat{\text{err}}_S(h) + \mathcal{E}_{u(h)}(m, \omega(u(h))) \cdot \delta$$

$$\leq \text{err}_D(h) + 2\mathcal{E}_{n(h)}(m, \omega(n(h)) \cdot \delta)$$

- Since $m \geq m_{VC}(\epsilon/2, \omega(n(h)) \cdot \delta, VC(H_n))$

$$\mathcal{E}_{n(h)}(m, \omega(n(h)) \cdot \delta) \leq \epsilon/2.$$

- So

$$\text{err}_D(\hat{h}) \leq \text{err}_D(h) + \epsilon.$$



Corollary:

H is non-uniformly learnable
if and only if H can be
expressed as a countable
union $H = \bigcup_{i=1}^{\infty} H_i$ where

H_1, H_2, \dots have finite VC-dim.

Minimum description length

- Suppose H is countable (or finite)
- Let $w: H \rightarrow [0,1]$ s.t. $\sum_{h \in H} w(h) = 1$
- For a fixed $h \in H$, Hoeffding's inequality states

$$|\text{err}_D(h) - \widehat{\text{err}}_S(h)| \leq \sqrt{\frac{\log(w(h)) + \log(2/\delta)}{2m}}$$
 w.p. $1 - w(h) \cdot \delta$
- By union bound over all $h \in H$,

$$\forall h \in H \quad |\text{err}_D(h) - \text{err}_S(h)| \leq \sqrt{\frac{\log(w(h)) + \log(1/\delta)}{2m}}$$

w.p. $1 - \delta$.

• SRM rule becomes

$$\hat{h} = \underset{h \in H}{\text{argmin}} \quad \widehat{\text{err}}_S(h) + \sqrt{\frac{\log(w(h)) + \log(1/\delta)}{2m}}$$

Prefix-Free codes

Definition: Function $c: H \rightarrow \{0,1\}^*$

is called a prefix-free code

iff $c(h)$ is NOT a prefix

of $c(h')$ for any $h, h' \in H$

$h \neq h'$.

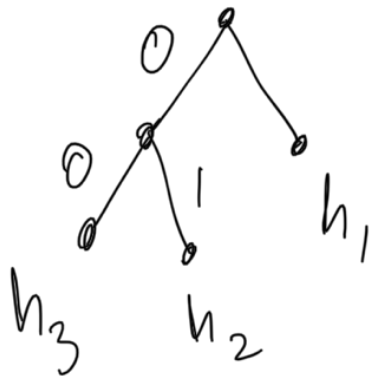
Example:

$$H = \{h_1, h_2, h_3\}$$

$$c(h_1) = 1$$

$$c(h_2) = 01$$

$$c(h_3) = 00$$



C is prefix tree

Example:

$$H = \{h_1, h_2, h_3\}$$

$$c(h_1) = 0$$

$$c(h_2) = 01$$

$$c(h_3) = 1$$

is NOT prefix-free

since $c(h_1)$ is prefix of $c(h_2)$.

Theorem (Kraft's inequality)

If $c: H \rightarrow \{0,1\}^*$ is prefix-free

then

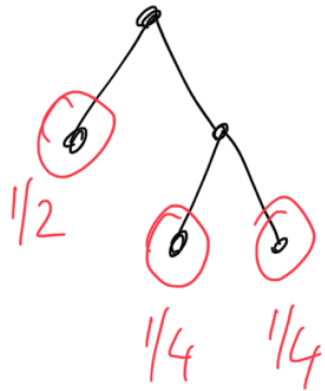
$$\sum_{h \in H} 2^{-|c(h)|} \leq 1.$$

Proof:

..... 1. 1. 1. T

- Fact: In a binary tree

$$\sum_{v \in \text{leaf}(T)} 2^{-\text{depth}(v)} \leq 1.$$



- Prefix-free code corresponds to leaves of a tree.
 - $\text{depth}(h) = |c(h)|$
-

- In MDL we can use

$$w(h) = 2^{-|c(h)|}$$

... is referred to as

- This is similar to ...

Occam's razor

- Occam in 13-th century suggested that

"simple explanations are better".