Introduction to Statistical Learning Theory Lecture 3

◆ロト ◆母 ▶ ◆ ヨ ▶ ◆ ヨ ● ● ● ●

Lecture 3

PAC learnability	VC dim: Examples	
0 0000 00	00 00	00 0000
Definition		

Reminder:

Definition (Growth function)

The growth function of \mathcal{H} , $\Pi_{\mathcal{H}}(m)$ is the size of the largest restriction of \mathcal{H} to a set of size m.

$$\Pi_{\mathcal{H}}(m) = \max\{|\mathcal{H}|_C|: C \subset \mathcal{X}, |C| = m\}$$

Theorem (Uniform convergence bound)

Let \mathcal{H} be a hypothesis set of $\{\pm 1\}$ valued functions and ℓ be the 0-1 loss, then for any distribution \mathcal{D} on $\mathcal{X} \times \{\pm 1\}$, any $\epsilon > 0$ and positive integer m, we have

$$P_{S \sim \mathcal{D}^m} \left(\exists h \in \mathcal{H} : \left| L_S(h) - L_{\mathcal{D}}(h) \right| \ge \epsilon \right) \le 4 \prod_{\mathcal{H}} (2m) \exp\left(-\frac{\epsilon^2 m}{8}\right)_{\mathcal{D} \subseteq \mathcal{A}}$$

PAC learnability
00
0000
00
Definition

In order to prove uniform convergence, and therefore PAC learnability, it is enough to show that the growth function is sub-exponential.

As we will see, the behavior $\Pi_{\mathcal{H}}(m)$ is greatly controlled by a single parameter - the VC dimension.

Definition (Shattering)

Let \mathcal{H} be a set of functions from \mathcal{X} to $\mathcal{Y} = \{\pm 1\}$. We say that \mathcal{H} shatters $C \subset \mathcal{X}$ if $\mathcal{H}|_C = 2^C$.

Definition (VC-dimension)

Let \mathcal{H} be a set of functions from \mathcal{X} to $\mathcal{Y} = \{\pm 1\}$. The VC-dimension of \mathcal{H} is the size of the largest finite set that \mathcal{H} shatters (or ∞ if there is no maximum).

So $VC(\mathcal{H}) = d \Leftrightarrow \Pi_{\mathcal{H}}(d) = 2^d \land \Pi_{\mathcal{H}}(d+1) < 2^{d+1} \Leftrightarrow \mathbb{R} \to \mathbb{R}$

PAC learnability	VC dim: Examples	
00 •000 00	00 00	00 0000
17701 1 11 11 11 11 11 11		

Theorem (Sauer–Shelah)

Let \mathcal{H} be a set of functions from \mathcal{X} to $\mathcal{Y} = \{\pm 1\}$ with VC-dimension $d < \infty$, then $\Pi_{\mathcal{H}}(m) \leq \sum_{k=1}^{d} {m \choose k}$

Notice that $S(m,d) = \sum_{k=0}^{d} {m \choose k}$ is the number of subset of size smaller or equal to d of a set of size m.

Proof.

We will show a stronger claim

$$|\mathcal{H}|_C| \le |\{B \subset C : \mathcal{H} \ shatters \ B \}| \le \sum_{k=0}^d {m \choose k}.$$

This is done by induction on m. For m = 1 the claim is trivial.

伺下 イヨト イヨト

PAC learnability	VC dim: Examples	
00 0000 00	00 00	00 0000
VC bounds on the growth function		

Proof (Cont.)

Let $C = \{x_1, ..., x_{m+1}\}$ and $\tilde{C} = \{x_1, ..., x_m\}$. Each function of $\mathcal{H}|_{\tilde{C}}$ corresponds to either one function in $\mathcal{H}|_C$ if it has a unique extension, or to two function if both extensions are possible.

Define $\mathcal{F} \subset \mathcal{H}|_{\tilde{C}}$ as all the function that correspond to two functions in $\mathcal{H}|_{C}$, then $|\mathcal{H}|_{C}| = |\mathcal{H}|_{\tilde{C}}| + |\mathcal{F}|$.

From our induction hypothesis $|\mathcal{H}|_{\tilde{C}}| \leq |\{B \subset \tilde{C} : \mathcal{H} \text{ shatters } B\}| = |\{B \subset C : \mathcal{H} \text{ shatters } B \land x_{m+1} \notin B\}|.$

For $\mathcal{F}: |\mathcal{F}| \leq |\{B \subset \tilde{C} : \mathcal{F} \text{ shatters } B\}|$. For each such B shattered by $\mathcal{F}, B \cup \{x_{m+1}\}$ is shattered by \mathcal{H} , so $|\mathcal{F}| \leq |\{B \subset C : \mathcal{H} \text{ shatters } B \land x_{m+1} \in B\}|$.

・ロッ ・雪ッ ・ヨッ

PAC learnability	VC dim: Examples	
00 00●0 00	00 00	00 0000
VC bounds on the growth function		

Our next step is bound $\Pi_{\mathcal{H}}(m)$ with a simpler bound:

Theorem

Let \mathcal{H} be a set of functions from \mathcal{X} to $\mathcal{Y} = \{\pm 1\}$ with *VC*-dimension $d < \infty$, then for m > d we have $\Pi_{\mathcal{H}}(m) \leq \left(\frac{em}{d}\right)^d$.

Proof.

We know from the Sauer-Shelah theorem that $\Pi_{\mathcal{H}}(m) \begin{cases} = 2^m & \text{if } m \leq d \\ \leq \sum_{k=0}^d \binom{m}{k} & \text{if } m > d \end{cases}$ We will prove $\sum_{k=0}^d \binom{m}{k} \leq \left(\frac{em}{d}\right)^d$ for $m \geq d$ to finish the proof.

・ロト ・四ト ・ヨト ・ヨト

э

PAC learnability	VC dim: Examples	
00 0000	00 00	00 0000
VC bounds on the growth function		

Proof (Cont.)

Need to prove
$$\sum_{k=0}^{d} {m \choose k} \le \left(\frac{em}{d}\right)^d$$
 for $m \ge d$.

$$\sum_{k=0}^{d} \binom{m}{k} \leq \left(\frac{m}{d}\right)^{d} \sum_{k=0}^{d} \binom{m}{k} \left(\frac{d}{m}\right)^{k} \leq \left(\frac{m}{d}\right)^{d} \sum_{k=0}^{m} \binom{m}{k} \left(\frac{d}{m}\right)^{k} = \\ \stackrel{*}{=} \left(\frac{m}{d}\right)^{d} \left(1 + \frac{d}{m}\right)^{m} \stackrel{**}{\leq} \left(\frac{em}{d}\right)^{d}$$

Where (*) is the binomial theorem, and (**) follows from Euler's inequality $\left(1 + \frac{1}{x}\right)^x < e$ for all x > 0.

We can combine all our results to show that if \mathcal{H} has VC-dimension $d < \infty$ then it is PAC learnable:

Theorem (PAC learnability of finite VC-dimension)

Let \mathcal{H} be a set of functions from \mathcal{X} to $\mathcal{Y} = \{\pm 1\}$ with VC-dimension $d < \infty$, then \mathcal{H} has the uniform convergence property with $\mathfrak{M}(\epsilon, \delta) = \mathcal{O}\left(\frac{d\ln(\frac{1}{\epsilon}) + \ln(\frac{1}{\delta})}{\epsilon^2}\right)$ and is therefore PAC learnable with the ERM algorithm.

Note: One can get (with some extra effort) a better bound without the $\ln(\frac{1}{\epsilon})$ factor.

PAC learnability	VC dim: Examples	
00 0000 0●	00 00	00 0000
PAC learnability of finite VC dimension		

Proof (sketch).

To prove uniform convergence we need to show that

$$P_{S \sim \mathcal{D}^m} \left(\exists h \in \mathcal{H} : \left| L_S(h) - L_{\mathcal{D}}(h) \right| \ge \epsilon \right) \le \delta \quad \forall m \ge \mathfrak{M}(\epsilon, \delta)$$

We already showed that

$$P_{S \sim \mathcal{D}^m} \left(\exists h \in \mathcal{H} : \left| L_S(h) - L_{\mathcal{D}}(h) \right| \ge \epsilon \right) \le 4 \Pi_{\mathcal{H}}(2m) \exp\left(-\frac{\epsilon^2 m}{8}\right)$$

<ロ> (日) (日) (日) (日) (日)

Using the inequality $\Pi_{\mathcal{H}}(m) \leq \left(\frac{em}{d}\right)^d$ and the inequality $\forall \alpha, x > 0 : \ln(x) \leq \alpha x - \ln(\alpha)$, we can show (with some algebra) that $4\Pi_{\mathcal{H}}(2m) \exp\left(-\frac{\epsilon^2 m}{8}\right) \leq \delta$.



Some examples: Example 1: $\mathcal{H} = 2^{\mathcal{X}}$ for infinite \mathcal{X} , $\Pi_{\mathcal{H}}(m) = 2^m$ Therefore $VC(\mathcal{H}) = \infty$.

Example 2: For finite \mathcal{H} , $\Pi_{\mathcal{H}}(m) \leq |\mathcal{H}| \Rightarrow VC(\mathcal{H}) \leq \log_2(|\mathcal{H}|)$.

Example 3: For $\mathcal{H} = \{h_a(x) = sign(x-a), a \in \mathbb{R}\},\$ $\Pi_{\mathcal{H}}(m) = m+1 \Rightarrow VC(\mathcal{H}) = 1.$

Example 4: For $\mathcal{H} = \{h_a^{\pm}(x) = sign(\pm x - a), a \in \mathbb{R}\},\$ $\Pi_{\mathcal{H}}(m) = 2m \Rightarrow VC(\mathcal{H}) = 2.$ Example 5: The class of axis aligned rectangles $h_{(x_1,x_2,y_1,y_2)}(x,y) := 1 \Leftrightarrow x_1 < x < x_2 \land y_1 < y < y_2$ for $x_1 < x_2, y_1 < y_2$.

It is easy to find a 4 element set that \mathcal{H} shatters. To show $VC(\mathcal{H}) = 4$, we need to show it cannot shatter any set of five elements. This can be done by observing that one point is always in the convex hall of the other points that cannot get zero if all others are one.

Example 6: The class of convex sets in the plan $h_C(x, y) := 1 \Leftrightarrow (x, y) \in C$ for a convex set $C \subset \mathbb{R}^2$. We can see that $VC(\mathcal{H}) = \infty$ by arranging points on the circle.



We will now compute the VC dimension on a simple, yet practical, hypothesis set - linear classifiers: $\mathcal{H}_d = \{h_w(x) = sign(\langle w, x \rangle) : w \in \mathbb{R}^d\}.$

To show that $VC(\mathcal{H}_d) = d$ it is enough to prove the following lemma:

Lemma

The vectors $x_1, ..., x_k \in \mathbb{R}^d$ are shattered by \mathcal{H}_d if and only if they are linearly independent.

Proof.

 \Rightarrow assume by contradiction that they are linearly dependent, so there exist some j such that $x_j = \sum_{i=1}^{j-1} \alpha_i x_i$. Any labeling y_i such that $\alpha_i y_i \ge 0$ has to have $y_j = 1$, therefore the set is not shattered - a contradiction.

 \Leftarrow Let X be the matrix with rows x_i^T , then the vector of labels given by any w is just $sign(X \cdot w)$. Our assumptions means that X has rank k and is therefore an onto mapping, and the set is shattered.

PAC learnability 00 0000 00	VC dim: Examples oo oo	No-Free-Lunch $\circ \circ$ $\circ \circ \circ \circ \circ$
Statement		

So far we have shown that if $VC(\mathcal{H}) < \infty$, then \mathcal{H} is learnable by the ERM algorithm. We will show bounds on what is possible to learn.

We start with the "No-Free-Lunch theorem" that shows there is no universal learner, i.e. any learning algorithm will fail on some distributions. This means that any (useful) learning algorithm has to have some assumptions on the task being learned.

Theorem (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:

1) There exists a function $f : \mathcal{X} \to \{0, 1\}$ such that $L_{\mathcal{D}}(f) = 0$. 2) With probability at least 1/7 over the choice of $S \sim \mathcal{D}^m$ we have that $L_{\mathcal{D}}(A(S)) \ge 1/8$. PAC learnability oo oocoo oo Proof VC dim: Examples $\begin{array}{c} \circ \circ \\ \circ \circ \end{array}$

No-Free-Lunch $^{\circ\circ}_{\bullet\circ\circ\circ}$

Proof.

First we can assume w.l.o.g. that $|\mathcal{X}| = 2m$ (simply by giving probability zero to excess points). There are $T = 2^{2m}$ functions $f : \mathcal{X} \to \{0, 1\}$, we will label them as $f_1, ..., f_T$. For each f_i define \mathcal{D}_i to be the distribution defined by f_i , i.e.

$$\mathcal{D}_i\left((x,y)\right) = \begin{cases} \frac{1}{2m} & if \ y = f_i(x) \\ 0 & otherwise. \end{cases}$$
(1)

The major step is to show that

 $\max_{i \in [T]} \mathbb{E}_{S \sim \mathcal{D}_i^m} [L_{\mathcal{D}_i}(A(S))] \geq 1/4.$ To prove this we will use the probabilistic method - we will show that by picking a distribution \mathcal{D}_i randomly, the expected loss is larger than 1/4.

We will need this final notation: There are $k = (2m)^m$ possible samples $S = (x_1, ..., x_m)$ which we label $S_1, ..., S_k$. For each distribution \mathcal{D}_i we define $S_i^i = ((x_1, f_i(x_1), ..., (x_m, f_i(x_m)))$.

▲□ ▶ ▲ □ ▶ ▲ □ ▶

э

Proof

Proof (cont.)

A simple changing of sums gives:

$$\mathbb{E}_{\mathcal{D}_i} \mathbb{E}_{S \sim \mathcal{D}_i^m} \left[L_{\mathcal{D}_i}(A(S)) \right] = \frac{1}{T} \sum_{i=1}^T \frac{1}{k} \sum_{j=1}^k L_{\mathcal{D}_i}(A(S_j^i))$$

We will finish the proof by showing $\frac{1}{T} \sum_i L_{\mathcal{D}_i}(A(S_j^i)) \ge 1/4$.

This can be done by looking at all the m unseen point $v_1, ..., v_m$ and to notice that for each v_p we can split the distributions into T/2 pairs f_{i_1}, f_{i_2} that agree on all the labels except on v_p .

▲聞を ▲ 国を ▲ 国を一

э

Proof

Proof (cont.)

A simple changing of sums gives:

$$\mathbb{E}_{\mathcal{D}_i} \mathbb{E}_{S \sim \mathcal{D}_i^m} \left[L_{\mathcal{D}_i}(A(S)) \right] = \frac{1}{T} \sum_{i=1}^T \frac{1}{k} \sum_{j=1}^k L_{\mathcal{D}_i}(A(S_j^i))$$
$$= \frac{1}{k} \sum_{j=1}^k \frac{1}{T} \sum_{i=1}^T L_{\mathcal{D}_i}(A(S_j^i)).$$

We will finish the proof by showing $\frac{1}{T} \sum_{i} L_{\mathcal{D}_i}(A(S_i^i)) \ge 1/4$.

This can be done by looking at all the m unseen point $v_1, ..., v_m$ and to notice that for each v_p we can split the distributions into T/2 pairs f_{i_1}, f_{i_2} that agree on all the labels except on v_p . VC dim: Examples $\begin{array}{c} \circ \circ \\ \circ \circ \end{array}$

No-Free-Lunch $\stackrel{\circ\circ}{_{\circ\circ}}_{\circ\circ\circ\circ}$

Proof

Proof (cont.)

As they agree on the sample, $S_j^{i_1} = S_j^{i_2}$ and the average loss between them is 1/2. Since half the points are unseen this gives a 1/4 loss bound. More rigorously:

$$\frac{1}{T} \sum_{i=1}^{T} L_{\mathcal{D}_i}(A(S_j^i)) \ge \frac{1}{T} \sum_{i=1}^{T} \frac{1}{2m} \sum_{p=1}^{m} \mathbb{1}[A(S_j^i)(v_p) = f_i(v_p)] = \frac{1}{T} \sum_{i=1}^{T/2} \frac{1}{2m} \sum_{p=1}^{m} \mathbb{1}[A(S_j^{i_1})(v_p) = f_{i_1}(v_p)] + \mathbb{1}[A(S_j^{i_2})(v_p) = f_{i_2}(v_p)] = \frac{1}{T} \sum_{i=1}^{T/2} \frac{1}{2m} \sum_{p=1}^{m} \mathbb{1} = \frac{1}{4}.$$

This shows that $\max_{i \in [T]} \mathbb{E}_{S \sim \mathcal{D}_i} [L_{\mathcal{D}_i}(A(S)] \ge 1/4.$

PAC learnability 00 0000 00	VC dim: Examples oo oo	No-Free-Lunch $\circ\circ$
Proof		

Proof (final).

We have $\max_{i \in [T]} \mathbb{E}_{S \sim \mathcal{D}_i} [L_{\mathcal{D}_i}(A(S)] \geq 1/4$. Mark that distribution as \mathcal{D} . To finish the proof we need to show that: with probability at least 1/7 over the choice of $S \sim \mathcal{D}^m$ we have that $L_{\mathcal{D}}(A(S)) \geq 1/8$. Using the law of total expectation we get:

$$\frac{1}{4} \leq \mathbb{E}_{S \sim \mathcal{D}} \left[L_{\mathcal{D}}(A(S)) \right] \leq P\left(L_{\mathcal{D}}(A(S) \geq \frac{1}{8}) \cdot 1 + \left(1 - P\left(L_{\mathcal{D}}(A(S) \geq \frac{1}{8}) \right) \cdot \frac{1}{8} \Rightarrow P\left(L_{\mathcal{D}}(A(S) \geq \frac{1}{8}) \geq \frac{1}{7} \right) \right)$$