

Lecture 5

- * HW 1 posted. Due in 3 weeks.
- * More details regarding presentations in \approx next week.
- * Scribing?
- * Last time : Concentration inequalities.
- * Today : VC dimension

Vapnik-Chervonenkis Dimension

Recap: For finite H ,

→ PAC learn with $O\left(\frac{\log(1/\delta)}{\varepsilon}\right)$ samples

→ agnostic-PAC learning with $O\left(\frac{\log(1/\delta/\varepsilon)}{\varepsilon^2}\right)$ samples

How about infinite H ?

→ Discretization: Linear classifier in \mathbb{R}^d

32 bit system $\rightarrow (2^{32})^d$ possible classifiers

→ VC dimension gives a nice way to handle ∞ classes.

Shattering

Def. (Restriction & Shattering)

The restriction of class H to a set of examples $C = \{c_1, \dots, c_n\} \in \mathcal{X}$ is a subset of $\{0,1\}^{|C|}$ given by $H_C = \{(h(c_1), \dots, h(c_n)) : h \in H\}$. We say that H shatters C if $|H_C| = 2^{|C|}$.

Shattering \Rightarrow all possible labellings are realized.

Corollary of No free lunch theorem.

Let \mathcal{H} be a hypothesis class & assume there exists a set $C \subseteq X$ of size $2n$ s.t. \mathcal{H} shatters C . Then if a distribution D over $X \times \{0,1\}$ & a predictor $h^* \in \mathcal{H}$ s.t. $R(h^*)=0$, but for any learning algorithm A , $\Pr_{S \sim D^n} [R(A(S)) \geq \frac{1}{8}] \geq \frac{1}{7}$

\mathcal{H} shatters a set of size $2n \Rightarrow$ cannot learn with n examples.

But if C is s.t. $|\mathcal{H}_C| \ll 2^{|C|}$ maybe can do something?

Idea: for any distribution supported on C , real hypothesis space under consideration is \mathcal{H}_C .

\mathcal{H}_C is finite!

If \mathcal{H}_C is small, then maybe can learn.

Def (VC dimension)

The VC dimension of a hypothesis class \mathcal{H} , denoted by $\text{VCdim}(\mathcal{H})$ is the size of the largest set $C \subseteq X$ that can be shattered by \mathcal{H} . If \mathcal{H} can shatter sets of arbitrary size, then $\text{VCdim}(\mathcal{H}) = \infty$.

Example

To show $\text{VCdim}(\mathcal{H}) = d$.

1. There exists some set C of size d that can be shattered by \mathcal{H} .
2. No set of size $d+1$ is shattered by \mathcal{H} .

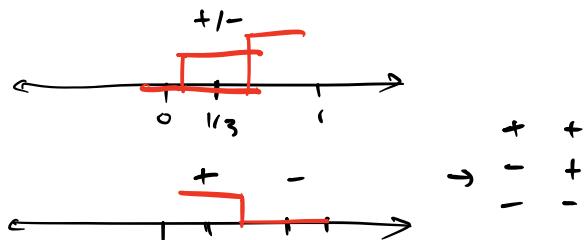
Example 1

Threshold functions:

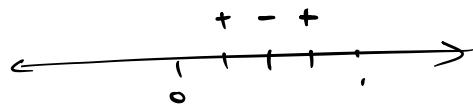
Let $X = [0, 1]$, $\mathcal{H} = \{h_b(+)=\mathbb{1}(x > b), b \in [0, 1]\}$

\mathcal{H} : set of thresholds in \mathbb{R} .

$$\text{VC dim}(\mathcal{H}) = 1$$



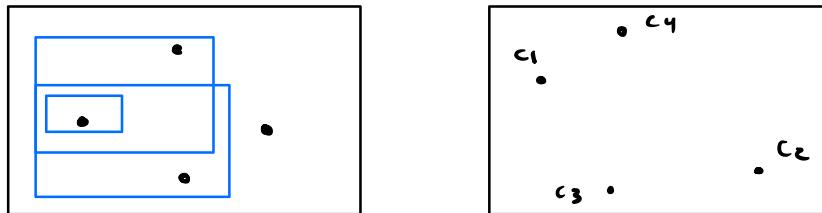
If we also allow reverse i.e. $\mathbb{1}(x < b)$, $\text{VCdim}(\mathcal{H}) = 2$.
(previous \mathcal{H} v/s new stuff's)



Example 2

Axis-aligned rectangles

$$X = \mathbb{R}^2, \quad \mathcal{H}_{a_1, a_2, b_1, b_2}(x_1, x_2) = \mathbf{1} \left(\begin{array}{l} a_1 \leq x_1 \leq b_1 \\ a_2 \leq x_2 \leq b_2 \end{array} \right)$$



$$\boxed{\begin{array}{l} \textcircled{1} \quad \left\{ \begin{array}{l} c_1: \min x_1 \\ c_2: \max x_1 \\ c_3: \min x_2 \\ c_4: \max x_2 \end{array} \right. \\ \textcircled{0} \quad \left\{ \begin{array}{l} c_5: \text{any other} \\ \text{(not unique minimizer)} \end{array} \right. \end{array}}$$

cannot realize this label.

Example 3

Finite classes.

For any finite hypothesis class \mathcal{H} , have $\text{VCdim}(\mathcal{H}) \leq \log(|\mathcal{H}|)$. This is because for any set C , $|\mathcal{H}_C| \leq |\mathcal{H}|$. Therefore, if $2^{|C|} > |\mathcal{H}|$, then we cannot shatter C .

VC Theorem

Let \mathcal{H} be a hypothesis class with $\text{VCdim}(\mathcal{H}) = d < \infty$.

Then there is a absolute constant $c > 0$ s.t.

\mathcal{H} has uniform convergence property with,

$$n_{\mathcal{H}}^{VC}(\epsilon, \delta) \leq c \cdot \frac{d \log(d/\epsilon) + \log(1/\delta)}{\epsilon^2}$$

Corollary

\mathcal{H} is agnostically-PAC learnable with

$$\mathcal{O}\left(\frac{d \log(d/\epsilon) + \log(1/\delta)}{\epsilon^2}\right) \text{ samples.}$$

Note

1) Possible to show $n_{\mathcal{H}}^{VC}(\epsilon, \delta) \leq c \cdot d + \frac{\log(1/\delta)}{\epsilon^2}$

2) This is for binary classification with 0/1 loss.
We'll see how to generalize to other losses /
beyond binary classification later.

for $d = \log(|\mathcal{H}|)$

$$n_{\mathcal{H}}^{VC}(\epsilon, \delta) \leq \mathcal{O}\left(\frac{\log(|\mathcal{H}|)}{\epsilon^2}\right) \text{ sample complexity}$$

Proof

- Outline:
- 1) for any set $C \subseteq X$, H_C : "Effective size" of H
 $|H_C| \approx |C|^d$.
 - 2) small "effective size" good for union bound to get VC.

Step 1: Polynomial growth of H_C .

Definition (Growth function)

The growth function of H , $T_H: \mathbb{N} \rightarrow \mathbb{N}$, is defined as

$$T_H(n) = \max_{C \subseteq X, |C|=n} |H_C|.$$

If $\text{VCdim}(H) = d$, then $T_H(n) = 2^n$ if $n \leq d$.

Sauer's lemma gives a good upper bound if $n \geq d$.

Lemma: (Sauer's Lemma)

$$\forall n, \text{VCdim}(H) = d, T_H(n) \leq \sum_{i=0}^d \binom{n}{i}.$$

For $n > d+1$, this implies $T_H(n) \leq \left(\frac{ne}{d}\right)^d$.

Proof: We instead show stronger inequality.

For any $C = \{c_1, \dots, c_n\}$ & any H .

$$|H_C| \leq |\{B \subseteq C : H \text{ shatters } B\}|. \quad (1)$$

This is sufficient since if $\text{VCdim}(H) = d$, H cannot shatter any set B of size $|B| > d$. There are $\binom{n}{d}$ subsets of size, we get our bound.

Prove (1) by induction.

Base case, $n=1$

Either 1) $|H_C| = 1$, LHS = RHS since one labelling, shatters $\{\emptyset\}$.

2) $|H_C| = 2$, two labellings shatters $\{\{\emptyset\}, \{c_1\}\}$

Induction step.

Assume that (1) holds for all sets of size $k < n$.

Let $C = \{c_1, \dots, c_n\}$ & $C' = \{c_2, \dots, c_n\}$

Define $\gamma_0 = \{(y_2, \dots, y_n) : (0, y_2, \dots, y_n) \in H_C \text{ or } (1, y_2, \dots, y_n) \in H_C\}$

$\gamma_1 = \{(y_2, \dots, y_n) : (0, y_2, \dots, y_n) \in H_{C'} \text{ and } (1, y_2, \dots, y_n) \in H_{C'}\}$

$$\text{Claim: } |\mathcal{H}_c| = |\mathcal{Y}_0| + |\mathcal{Y}_1|$$

True because (y_2, \dots, y_n) is counted once in \mathcal{Y}_0 , but again in \mathcal{Y}_1 if can be shattered.

By induction we get,

$$\begin{aligned} |\mathcal{Y}_0| &\leq |\{B \subseteq C : \mathcal{H} \text{ shatters } B\}| \\ &= |\{B \subseteq C : c_i \notin B \text{ AND } \mathcal{H} \text{ shatters } B\}|. \end{aligned}$$

for \mathcal{Y}_1 , define $\mathcal{H}' \subseteq \mathcal{H}$ to be

$$\begin{aligned} \mathcal{H}' &= \{h \in \mathcal{H}, \exists h' \in \mathcal{H} \text{ s.t. } ((1-h'(c_1), h'(c_2), \dots, h'(c_n)) \\ &\quad = (h(c_1), h(c_2), \dots, h(c_n)))\} \end{aligned}$$

\mathcal{H}' : all hypothesis for which the hypothesis that agrees everywhere in C except c_i is also in \mathcal{H} .

Note : 1) IF \mathcal{H}' shatters $B \subseteq C'$ then it also shatters $B \cup \{c_i\}$.

$$2) \mathcal{Y}_1 = \mathcal{H}'_{c_i}$$

$$\begin{aligned} |\mathcal{Y}_1| &= |\mathcal{H}'_{c_i}| \leq |\{B \subseteq C' : \mathcal{H}' \text{ shatters } B\}| \quad (\text{By (1)}) \\ &= |\{B \subseteq C' : \mathcal{H}' \text{ shatters } B \cup \{c_i\}\}| \\ &= |\{B \subseteq C : c_i \in B \text{ AND } \mathcal{H} \text{ shatters } B\}| \end{aligned}$$

$$\leq |\{B \subseteq C : c_i \in B \text{ AND } \mathcal{H} \text{ shatters } B\}|$$

$$\therefore |\mathcal{H}_c| = |\mathcal{Y}_0| + |\mathcal{Y}_1|$$

$$\leq |\{B \subseteq C : c_i \notin B \text{ AND } \mathcal{H} \text{ shatters } B\}|$$

$$+ |\{B \subseteq C : c_i \in B \text{ AND } \mathcal{H} \text{ shatters } B\}|$$

$$= |\{B \subseteq C : \mathcal{H} \text{ shatters } B\}|,$$

which completes our proof. \blacksquare

Step 2: Symmetrization

Lemma: For a class \mathcal{H} with growth functions $T_{\mathcal{H}}$,

$$\mathbb{E}_{S \sim D^n} \sup_{h \in \mathcal{H}} |R(h) - \hat{R}_S(h)| \leq \sqrt{\frac{2 \log(2T_{\mathcal{H}}(2n))}{n}}.$$

Once we have this,

we can Markov's inequality to get a high probability statement.

$$\Pr \left[\sup_{h \in \mathcal{H}} |R(h) - \hat{R}_S(h)| > t \right] \leq \frac{\sqrt{\frac{2 \log(2T_{\mathcal{H}}(2n))}{n}}}{t}.$$

But we'll use McDiarmid's to get something better.