

Lecture 6

- * Please sign up to scribe if you haven't yet!
- * Sign up for presentation topics tomorrow
- * Last time
 - VC Dimension
 - VC Theorem

Today :

- finish proof of VC Theorem
- Rademacher complexity

VC Theorem

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

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

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

Proof

Outline: 1) for any set $C \subseteq \mathcal{X}$, \mathcal{H}_C : "Effective size" of \mathcal{H}
 $|\mathcal{H}_C| \approx |\mathcal{C}|^d$.

2) small "effective size" good for union bound to get VC.

Def. (Restriction)

The restriction of class \mathcal{H} to a set of examples $C = \{c_1, \dots, c_n\} \subseteq \mathcal{X}$ is a subset of $\{0,1\}^{|C|}$ given by $\mathcal{H}_C = \{(h(c_1), \dots, h(c_n)) : h \in \mathcal{H}\}$

Step 1: Polynomial growth of \mathcal{H}_c .

Definition (growth function)

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

$$T_{\mathcal{H}}(n) = \max_{C \subseteq \mathcal{X}, |C| \leq n} |\mathcal{H}_C|.$$

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

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

Lemma : (Sauer's Lemma)

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

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

Step 2: Symmetrization

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

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

Proof

We will use the idea of symmetrization.

Let $S' = \{(+i', g_i'), i' \in [n]\}$ be identically distributed as S .

Note that $\mathbb{E}_{S'} [\hat{R}_{S'}(h)] = R(h)$.

Therefore

$$\mathbb{E}_S \sup_{h \in \mathcal{H}} |R(h) - \hat{R}_S(h)| = \mathbb{E}_S \sup_{h \in \mathcal{H}} \left| \mathbb{E}_{S'} [\hat{R}_{S'}(h)] - \hat{R}_S(h) \right|$$

Fix S .

$$\text{Claim: } \sup_{h \in \mathcal{H}} \left| \mathbb{E}_{S'} [\hat{R}_{S'}(h)] \right| \leq \mathbb{E}_{S'} \sup_{h \in \mathcal{H}} |\hat{R}_{S'}(h)|.$$

Proof.

1. $|\cdot|$ is a convex function

& sup/max of convex functions is convex

$\therefore \sup_{h \in \mathcal{H}} \left| \mathbb{E}_{S'} [\hat{R}_{S'}(h)] \right|$ is a convex function of $\hat{R}_{S'}(h)$

By Jensen's claim follows.

(Jensen's,
 $f(\mathbb{E}(x)) \leq \mathbb{E}(f(x))$
if f convex)

$$\begin{aligned} \therefore \mathbb{E}_S \sup_{h \in \mathcal{H}} |R(h) - \hat{R}_S(h)| &\leq \mathbb{E}_{S, S'} \sup_{h \in \mathcal{H}} |\hat{R}_{S'}(h) - \hat{R}_S(h)| \\ &= \mathbb{E}_{S, S'} \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \left(\mathbb{1}_{\{h(x_i) \neq y_i'\}} - \mathbb{1}_{\{h(x_i) \neq y_i\}} \right) \right|. \end{aligned}$$

independent

Let $\sigma_{1:n} = \{\sigma_1, \dots, \sigma_n\}$ be Rademacher r.v. i.e. $\text{Unif}(\{-1\})$.

Since (x_i, y_i) & (x'_i, y'_i) are i.i.d.

$$\begin{aligned} \therefore \mathbb{1}_{\{h(x'_i) \neq y'_i\}} - \mathbb{1}_{\{h(x_i) \neq y_i\}} \\ \sim \mathbb{1}_{\{h(x_i) \neq y_i\}} - \mathbb{1}_{\{h(x'_i) \neq y'_i\}}. \end{aligned}$$

$$\begin{aligned} \therefore \mathbb{E}_S \sup_{h \in \mathcal{H}} |R(h) - \hat{R}_S(h)| \\ \leq \mathbb{E}_{\sigma_{1:n}} \mathbb{E}_{S, S'} \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \left(\mathbb{1}_{\{h(x_i) \neq y_i'\}} - \mathbb{1}_{\{h(x_i) \neq y_i\}} \right) \right| \\ = \mathbb{E}_{S, S'} \mathbb{E}_{\sigma_{1:n}} \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \left(\mathbb{1}_{\{h(x_i) \neq y_i'\}} - \mathbb{1}_{\{h(x_i) \neq y_i\}} \right) \right| \end{aligned}$$

Now fix s, s' & let C be the set of examples appearing in s, s' (the union). Note $|C| \leq 2n$.

Key idea We can replace sup over \mathcal{H} by max over \mathcal{H}_C .

$$\begin{aligned} & \mathbb{E}_{\sigma_{i:n}} \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \left(\mathbb{1}_{\{h(x_i) \neq y_i\}} - \mathbb{1}_{\{h(x_i) = y_i\}} \right) \right| \\ &= \mathbb{E}_{\sigma_{i:n}} \max_{h \in \mathcal{H}_C} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \left(\mathbb{1}_{\{h(x_i) \neq y_i\}} - \mathbb{1}_{\{h(x_i) = y_i\}} \right) \right| \end{aligned}$$

Let $\theta_h = \frac{1}{n} \sum_{i=1}^n \sigma_i \left(\mathbb{1}_{\{h(x_i) \neq y_i\}} - \mathbb{1}_{\{h(x_i) = y_i\}} \right)$

$$\mathbb{E}_s \sup_{h \in \mathcal{H}} |R(h) - \hat{R}_s(h)| \leq \mathbb{E}_{s,s'} \mathbb{E}_{\sigma_{i:n}} \max_{h \in \mathcal{H}_C} |\theta_h|$$

Want to bound $\mathbb{E}_{\sigma_{i:n}} \max_{h \in \mathcal{H}_C} |\theta_h|$

Lemma : (max of sub-Gaussian). If x_1, \dots, x_m are mean 0

σ -sub-Gaussian (not necessarily independent), then

$$\mathbb{E} \max_i x_i \leq \sigma \sqrt{2 \log(m)}$$

Proof

$$\begin{aligned}\mathbb{E} \max_i x_i &= \frac{1}{\lambda} \log \exp \left(\lambda \mathbb{E} [\max_i x_i] \right) + \lambda \\ &\leq \frac{1}{\lambda} \log \mathbb{E} [\exp (\lambda \max_i x_i)] \quad (\text{Jensen's}) \\ &\leq \frac{1}{\lambda} \log \mathbb{E} \left[\sum_{i=1}^m \exp (\lambda x_i) \right] \\ &= \frac{1}{\lambda} \log \left(\sum_{i=1}^m \mathbb{E} (\exp (\lambda x_i)) \right) \\ &\leq \frac{1}{\lambda} \log \left(\sum_{i=1}^m \exp \left(\frac{\lambda^2 \sigma^2}{2} \right) \right)\end{aligned}$$

$$\lambda = \frac{\sqrt{2 \log(m)}}{\sigma}$$

$$\mathbb{E} \max_i x_i \leq \frac{\sigma}{\sqrt{2 \log(m)}} \log \left(\sum_{i=1}^m \exp (\log m) \right)$$

$$= \sigma \cdot \frac{2 \log m}{\sqrt{2 \log m}} = \sigma \sqrt{2 \log m}.$$

•

Claim: θ_n is sub-Gaussian with parameter $\frac{1}{\sqrt{n}}$, $\mathbb{E}[\theta_n] = 0$.

Proof

$$\theta_n = \sum_{i=1}^n \frac{\epsilon_i}{n} \left(\mathbb{1}\{h(x_i) \neq y_i\} - \mathbb{1}\{h(x_i) = y_i\} \right)$$

$$\mathbb{E}[\theta_n] = 0. \quad (\sigma_i \sim \text{Unif}(\{-1, 1\}))$$

$$\theta_n = \sum_{i=1}^n \underbrace{\frac{\epsilon_i}{n} \left(\mathbb{1}\{h(x_i) \neq y_i\} - \mathbb{1}\{h(x_i) = y_i\} \right)}$$

each term is sub-Gaussian with parameter $\frac{1}{n}$.

$\therefore \theta_n$ is sub-Gaussian with parameter $\left(\sum_{i=1}^n \frac{1}{n^2} \right)^{1/2}$

$$= \frac{1}{\sqrt{n}}.$$

(θ_n is sum of independent sub-Gaussians with parameter $\frac{1}{n}$).

Claim: $\mathbb{E}_{\sigma_{i:n} \text{ mat } h \in \mathcal{H}_C} |\theta_n| \leq \frac{1}{\sqrt{n}} \sqrt{2 \log(2|\mathcal{H}_C|)}$

Proof

$$\mathbb{E}_{\sigma_{i:n} \text{ mat } h \in \mathcal{H}_C} |\theta_n| = \mathbb{E}_{\sigma_{i:n} \text{ mat } h \in \mathcal{H}_C} \mathbb{E}_{\theta_n, -\theta_n} \{ \theta_n \}$$

Recall that if θ_n is sub-Gaussian

$\Rightarrow -\theta_n$ is sub-Gaussian (with same parameter)

$$\therefore \mathbb{E}_{\sigma_{1:n}} \max_{h \in \mathcal{H}_c} |\theta_h| \leq \frac{1}{\sqrt{n}} \sqrt{2 \log(2|\mathcal{H}_c|)}$$

•

We have shown

$$\begin{aligned} \mathbb{E}_{S, \sigma_1^n} [\mathbb{E}_{\sigma_{1:n}} \max_{h \in \mathcal{H}_c} |\theta_h|] &\leq \sqrt{\frac{2 \log(2|\mathcal{H}_c|)}{n}} \\ \therefore \mathbb{E}_S \sup_{h \in \mathcal{H}_c} |R(h) - \hat{R}_S(h)| &\leq \sqrt{\frac{2 \log(2|\mathcal{H}_c|)}{n}} \end{aligned}$$

Note that $|\mathcal{H}_c| \leq T_H(2n)$

since $|C|=2n$.

Finishes proof by plugging in bound.

•

Step 3. McDiarmids.

$$\text{Define } f(s) = \sup_{h \in \mathcal{H}} |R(h) - \hat{R}_S(h)|$$

Observe that $f(s)$ satisfies the bounded difference property with constant $\frac{1}{n}$.

(since changing (x_i, y_i) can only change $\hat{R}_S(h)$ by $\frac{1}{n}$ for any $h \in \mathcal{H}$ \therefore the max also changes by at most $\frac{1}{n}$)

Using McDiarmids

$$\Pr[f(s) - \mathbb{E}[f(s)] > t] \leq 2e^{-t^2/2n}$$

Choose $t = \sqrt{\frac{\log(2/\delta)}{2n}}$ to get failure probability δ .

\therefore with prob $1-\delta$

$$\sup_{h \in H} |R(h) - \hat{R}_S(h)| \leq \sqrt{\frac{2\log(2T_H(2n))}{n}} + \sqrt{\frac{\log(2/\delta)}{2n}}$$

Step 4. finish proof.

Using Sauer's lemma, for $n \geq d+1$, $T_H(n) \leq \left(\frac{ne}{d}\right)^d$.

Plugging this in, with probability $(1-\delta)$,

$$\sup_{h \in H} |R(h) - \hat{R}_S(h)| \leq \sqrt{\frac{2d \log(2nd)}{n}} + \sqrt{\frac{\log(2/\delta)}{2n}}.$$

\therefore for $n \geq \mathcal{O}\left(d \frac{\log(d/\varepsilon) + \log(1/\delta)}{\varepsilon^2}\right)$ the RHS $\leq \varepsilon$.

[Exercise: finish this]

Rademacher complexity

Let us recall the proof of the VC theorem.

We wanted to bound

$$\mathbb{E}_s \sup_{h \in \mathcal{H}} |R(h) - \hat{R}_s(h)|.$$

This quantity is called an empirical process

Empirical process theory studies such quantities,

$$\begin{aligned} \mathbb{E}_s \sup_{h \in \mathcal{H}} (R(h) - \hat{R}_s(h)) &\leq \mathbb{E}_{s, s'} \sup_{h \in \mathcal{H}} (\hat{R}_{s'}(h) - \hat{R}_s(h)) \\ &= \mathbb{E}_{s, s'} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \left(\mathbb{1}\{h(x_i) \neq y_i\} - \mathbb{1}\{h(x_i) = y_i\} \right) \\ &= \mathbb{E}_{\sigma_{1:n}} \mathbb{E}_{s, s'} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_s \left(\mathbb{1}\{h(x_i) \neq y_i\} - \mathbb{1}\{h(x_i) = y_i\} \right) \\ &\leq \mathbb{E}_s \mathbb{E}_{\sigma_{1:n}} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbb{1}\{h(x_i) \neq y_i\} \\ &\quad + \mathbb{E}_s \mathbb{E}_{\sigma_{1:n}} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (\sigma_i) \mathbb{1}\{h(x_i) \neq y_i\} \end{aligned}$$

$$\mathbb{E}_s \sup_{h \in \mathcal{H}} R(h) - \hat{R}_s(h) \leq 2 \mathbb{E}_s \mathbb{E}_{\sigma_{1:n}} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbb{1}\{h(x_i) \neq y_i\}$$

Let $Z = X + Y$

F : function class $Z \rightarrow \mathbb{R}$

D : distribution over Z .

Definition (Rademacher Complexity).

Let F be a family of real-valued functions

$f: Z \rightarrow \mathbb{R}$ where $Z = X + Y$. Then the Rademacher complexity $R(F)$ is defined as

$$R(F) = \frac{1}{n} \mathbb{E}_{\sigma \sim \{\pm 1\}^n} \left[\sup_{f \in F} \sum_{i=1}^n \sigma_i f(z_i) \right].$$

More generally, given a set of vectors $A \subset \mathbb{R}^n$, the Rademacher complexity $R(A)$ is defined as

$$R(A) = \frac{1}{n} \mathbb{E}_{\sigma \sim \{\pm 1\}^n} \sup_{a \in A} \sum_{i=1}^n \sigma_i a_i.$$

Intuition: $R(F)$ captures how well can F fit random noise.

If F can fit random noise,
 F will probably overfit.