

Lecture 23

Last time

- ▷ Vapnik-K-Chervonenkis (VC) Theory.

Today

- ▷ Controlling the VC dimension
- ▷ Generalization for linear models.

Controlling the VC dimension

Let G be a class of real valued functions $\{g: X \rightarrow \mathbb{R}\}$. Define

$$S_g = \{x \in X \mid g(x) \leq 0\}$$

$$S(G) = \{S_g \mid g \in G\}.$$

Proposition: Let G be a vectors space of functions $g: \mathbb{R}^d \rightarrow \mathbb{R}$ with $\dim(G) < \infty$. Then,

$$\text{VC}(S(G)) \leq \dim(G).$$

Proof: Set $n = \dim(G) + 1$ and suppose $S(G)$ shatters some $X = \{x_1, \dots, x_n\}$.

Define

$$L: G \rightarrow \mathbb{R}^n$$
$$L(g) \mapsto \begin{pmatrix} g(x_1) \\ \vdots \\ g(x_n) \end{pmatrix}.$$

Since $n > \dim(G)$, there is
 $\gamma \neq 0 \in \mathbb{R}^n$ s.t. $\langle \gamma, L(g) \rangle = 0$
for all $g \in G$. Therefore,

$$-\sum_{\{i : \gamma_i \leq 0\}} \gamma_i g(x_i) = \sum_{\{i : \gamma_i > 0\}} \gamma_i g(x_i) \quad \forall g \in G.$$

Δ_1 Δ_2 T_1 T_2

WLOG assume $\gamma_i > 0$ for some i .
Since $S(Lg)$ shatters X , there is a
function g s.t.

$$\begin{aligned} g(x_i) \leq 0 & \quad \forall i \in \Delta_1 \quad \text{and} \\ g(x_i) > 0 & \quad \forall i \in \Delta_2 \end{aligned}$$

so

$$0 \geq T_1 = T_2 > 0.$$

Q

Example (Half spaces)

Define

$$S_{a,b} = \{x \in \mathbb{R}^d \mid \langle a, x \rangle + b \leq 0\}$$

$$S = \{S_{a,b} \mid a \in \mathbb{R}^d, b \in \mathbb{R}\}$$

$$\mathcal{G} = \{x \mapsto \langle a, x \rangle + b\}.$$

Then, \leftarrow HWS

$$VC(S) \leq \dim(\mathcal{G}) = d+1.$$

→

Example (Balls):

Define

$$S_{a,b} = \{x \in \mathbb{R}^d \mid \|x - a\|_2 \leq b\}$$

$$S = \{S_{a,b} \mid a \in \mathbb{R}^d, b \geq 0\}.$$

Consider functions

$$g_{a,b}(x) = \|x - a\|^2 - b^2$$

$$= \|x\|^2 + 2\langle a, x \rangle + \|a\|^2 - b^2.$$

This is not a subspace. But it is contained in one.

The trick is to define

$$g_c(x) = \langle c, \theta(x) \rangle \text{ with } c \in \mathbb{R}^{d+2}$$

and

$$\theta(x) = (1, x_1, \dots, x_d, \|x\|_2^2).$$

Then, it is clear that

$$\{g_{a,b} \mid \begin{matrix} a \in \mathbb{R}^d \\ b \in \mathbb{R} \end{matrix}\} \subseteq \{g_c \mid c \in \mathbb{R}^{d+1}\}$$

↑ subspace of
dim $d+2$

So

$$VC(S) \leq d+2.$$

The correct value is $d+1$ (much harder to prove)

As we have seen with these simple examples, we often get guarantees of the form

$$E \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum f(x_i) - E f(x) \right| \right] = O \left(\sqrt{\frac{d \log(n)}{n}} \right)$$

where d is the input dimension. A natural question is

Can we get dimension independent

bounds for interesting \mathcal{F} ?

Dimension free bounds

We will derive bounds for two types of losses: (1) linear models and (2) strongly convex functions.

Linear models

Consider the problem $\min_{\theta \in \Theta} \mathbb{E}_{(a,b) \sim P} l(\langle \theta, a \rangle, b)$ l is 1-Lipschitz in its first coordinate

Thus, we aim to prove generalization bounds by controlling $R_n(\mathcal{F})$ of $\mathcal{F} = \{ (a, b) \mapsto l(\langle \theta, a \rangle, b) \mid \theta \in \Theta \}$.

Recall that we had

$$\begin{aligned} \sup_{\theta \in \Theta} \left| \frac{1}{n} \sum l(\langle \theta, a_i \rangle, b_i) - \mathbb{E} l(\langle \theta, a \rangle, b) \right| \\ \leq 2 \mathbb{E}_{a_i, b_i \sim P} \sup \left| \sum_{i=1}^n \varepsilon_i l(\langle \theta, a_i \rangle, b_i) \right|. \end{aligned}$$

In turn we can also prove the same thing without absolute values, which

↑
Prove it

still gives us interesting information. (Why?)
 Let \bar{R} and R_n be Rademacher complexities without absolute values.

Theorem (Contraction principle): Consider a set $A \subseteq \mathbb{R}^n$ and let $\varphi_i : \mathbb{R} \rightarrow \mathbb{R}$ be ρ -Lipschitz functions and let

$$A' = \{(\varphi_1(a_1), \dots, \varphi_n(a_n)) \mid a \in A\}.$$

Then,

$$\bar{R}(A') = \mathbb{E} \sup_{a \in A'} \sum \varepsilon_i a_i \leq \mathbb{E} \sup_{a \in A} \sum \varepsilon_i a_i'' \quad \text{---} \\ \text{Notice it is } A'.$$

Notice that this theorem allows us to take $\varphi_i(a) = l(\langle \theta, a \rangle, b_i)$ and

$$\bar{R}_n(F) \leq \bar{R}_n(\{a \mapsto \langle \theta, a \rangle \mid \theta \in \Theta\}).$$

Proof of the contraction principle

WLOG $\rho=1$. For any $\psi : \mathbb{R}^n \rightarrow \mathbb{R}^n$ with ψ_i 1-Lipschitz define

$$\Psi(A) = \{(\psi_1(a_1), \dots, \psi_n(a_n)) \mid a \in A\}.$$

Let $\Psi = (\varphi_1, \dots, \varphi_n)$, and $\tilde{\Psi} = (I, \varphi_2, \dots,$

Let $\Psi = (\varphi_1, \dots, \varphi_n)$, and $\tilde{\Psi} = (\mathbb{I}, \varphi_2, \dots, \varphi_n)$. By the permutation invariance of \bar{R} , it suffices to prove

$$(\star) \quad \bar{R}(\Psi(A)) \leq \bar{R}(\tilde{\Psi}(A)).$$

Expanding

$$\bar{R}(\Psi(A)) = \mathbb{E}_\varepsilon \sup_{a \in A} \sum_{i=1}^n \varepsilon_i \varphi_i(a_i)$$

$$\begin{aligned} \text{Def } \rightarrow &= \frac{1}{2} \left[\mathbb{E}_{\varepsilon_2:n} \sup_{a \in A} \varphi_1(a_1) + \sum_{i=2}^n \varepsilon_i \varphi_i(a_i) \right. \\ &\quad \left. + \mathbb{E}_{\varepsilon_2:n} \sup_{a \in A} -\varphi_1(a_1) + \sum_{i=2}^n \varepsilon_i \varphi_i(a_i) \right] \end{aligned}$$

$$\begin{aligned} &= \frac{1}{2} \left[\mathbb{E}_{\varepsilon_2:n} \sup_{a, \tilde{a} \in A} \varphi_1(a_1) - \varphi_1(\tilde{a}_1) \right. \\ &\quad \left. + \sum_{i=2}^n \varepsilon_i (\varphi_i(a_i) + \varphi_i(\tilde{a}_i)) \right] \end{aligned}$$

1-Lipschitz

$$\leq \frac{1}{2} \left[\mathbb{E}_{\varepsilon_2:n} \sup |a_1 - \tilde{a}_1| \right. \\ \left. + \sum_{i=2}^n \varepsilon_i (\varphi_i(a_i) + \varphi_i(\tilde{a}_i)) \right]$$

Because of

the \sup we can

ensure $a_1 - \tilde{a}_1 \geq 0$.

$$\begin{aligned} &= \frac{1}{2} \left[\mathbb{E}_{\varepsilon_2:n} \sup |a_1 - \tilde{a}_1| \right. \\ &\quad \left. + \sum_{i=2}^n \varepsilon_i (\varphi_i(a_i) + \varphi_i(\tilde{a}_i)) \right] \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2} \left[\mathbb{E} \sup_{\varepsilon_2 \in \mathbb{N}} a_1 + \sum_{i=2}^n \varepsilon_i \Psi_i(a_i) \right. \\
&\quad \left. + \mathbb{E} \sup_{\varepsilon_2 \in \mathbb{N}} -a_1 + \sum_{i=2}^n \varepsilon_i \Psi_i(a_i) \right] \\
\text{Def } \rightarrow &= \bar{R}(\hat{\Psi}(A)).
\end{aligned}$$

□