

Quantity of interest

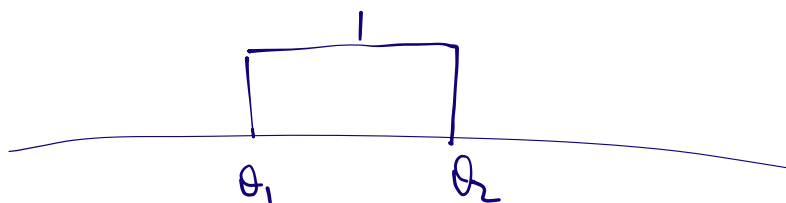
$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E} f(X) \right|$$

• Symmetrization (Rademacher complexity)

• Dudley chaining $|f(x)| \leq F(x) \quad \forall f$

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \dots \right] \leq c \sqrt{\frac{\mathbb{E} [F(X)^2]}{n}} \int_0^1 \sqrt{\log \sup_{\mathcal{Q}} N(\epsilon F, L^2(\mathcal{Q}); \mathcal{F}, L^2(\mathcal{Q}))} d\epsilon$$

eg. $\mathcal{F} := \{ \mathbb{1}_{\{x \in [\theta_1, \theta_2]\}} : \theta_1, \theta_2 \in \mathbb{R} \}$
 covering in (θ_1, θ_2)



\downarrow

covering in \mathcal{F} .

Vapnik - Chernovenkis (Binary functions $X \rightarrow \{0, 1\}$)

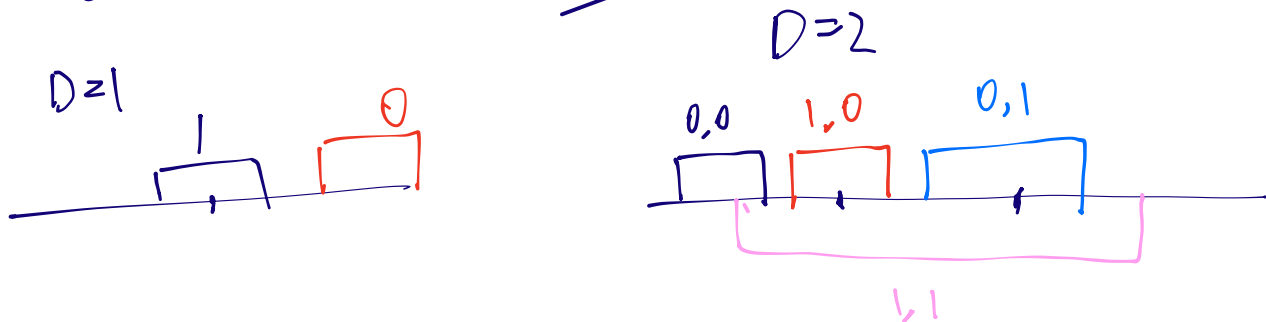
Def — We say (x_1, x_2, \dots, x_n) is shattered

by \mathcal{F} , if

$$\left\{ (f(x_1), f(x_2), \dots, f(x_n)) \mid f \in \mathcal{F} \right\} = \{0, 1\}^n.$$

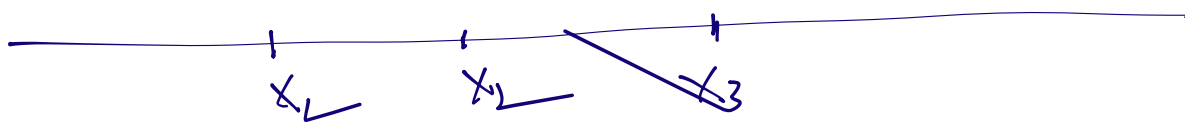
Def. $VC(\mathcal{F}) :=$ largest value of D st.
 $\exists x_1, x_2, \dots, x_D \in X$ shattered by \mathcal{F} .

eg. indicator of line segment



$(1, 0, 1)$ is not achievable.

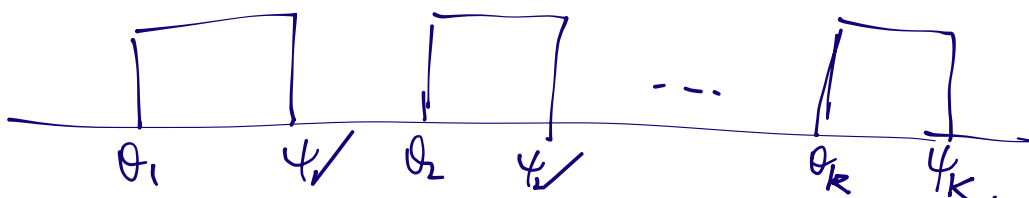
$D=3$



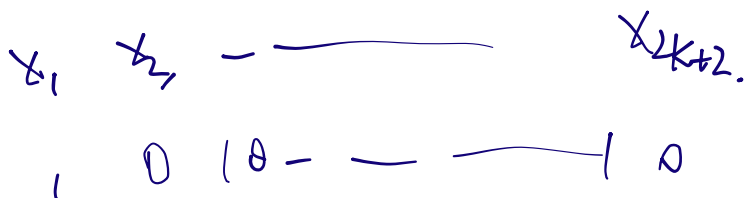
$$VC(\mathcal{F}) = 2.$$

eg. $\mathcal{F} := \left\{ \mathbb{I}_{x \in \bigcup_{i=1}^K [a_i, b_i]} \right\}$.

$2K$ (label changes)



$2k+2$ polynomials



$$VC(\mathcal{F}) < 2k+2.$$

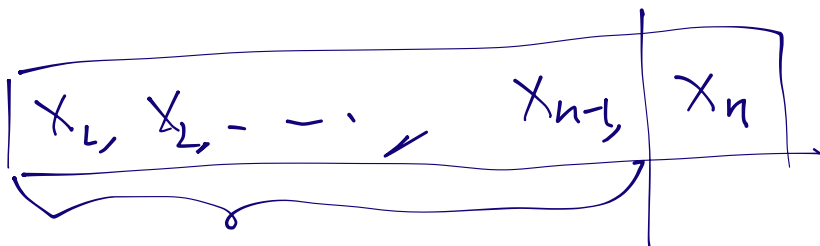
Sauer's lemma If $VC(\mathcal{F}) \leq D$ (for $n \geq D$).

Then $\left| \left\{ (f(x_1), \dots, f(x_n)) \in \mathcal{F} \right\} \right| \leq \sum_{i=0}^D \binom{n}{i} \leq \left(\frac{en}{D} \right)^D$

Proof: Induction $n \geq D$.

$n=1$ trivial

Suppose the conclusion holds for $(n-1, D-1)$ and (n, D) .



Consider binary sequence b_1, b_2, \dots, b_m

- b_1, \dots, b_m is not achieved by $\mathcal{F}(x_1, \dots, x_n)$
- ① $b_1, \dots, b_{m-D}, b_1, \dots, b_{m-1} \in \mathcal{F}(x_1, \dots, x_n)$ *ignore it*
 - ② Only one of them

Construct two sub classes $\mathcal{F}_1, \mathcal{F}_2$ as follows:

① $f_1, f_2 \in \mathcal{F}$ $f_1(x_1, \dots, x_n) = b_1, b_2, \dots, b_{n-1}$ \circ
 $f_2(x_1, \dots, x_n) = b_1, \dots, b_{n-1}$

Add f_1 to \mathcal{F}_1

Add f_2 to \mathcal{F}_2 .

② $f(x_1, \dots, x_n) = b_1, \dots, b_{n-1}, *$
 Add f to \mathcal{F}_1 .

We regard \mathcal{F}_1 and \mathcal{F}_2 as function classes with domain (x_1, x_2, \dots, x_n) .

• $|\mathcal{F}(x_1, \dots, x_n)| = |\mathcal{F}_1| + |\mathcal{F}_2|$.

• $|\mathcal{F}_1| = |\mathcal{F}(x_1, \dots, x_{n-1})| \leq \sum_{n-1}(D)$.

• $\forall f \in \mathcal{F}_2, \exists \tilde{f} \in \mathcal{F}_1$ s.t. f and \tilde{f} disagree only on x_n .

Claim: $VC(\mathcal{F}_2) \leq D-1$. \leftarrow

Suppose $(x_{k_1}, x_{k_2}, \dots, x_{k_D})$ shattered by \mathcal{F}_2 .
 $(x_{k_1}, x_{k_2}, \dots, x_{k_D}, x_n)$ shattered by \mathcal{F}_1 ,
 $(D+1)$ length.

$$|\mathcal{F}_2| \leq \bar{\Phi}_{n-1}(D_2).$$

$$\text{So } |\mathcal{F}(x_1, \dots, x_n)| \leq \bar{\Phi}_n(D_1) + \bar{\Phi}_{n-1}(D_2).$$

eg. $VC(\mathcal{F}_1) \leq D_1$ $VC(\mathcal{F}_2) \leq D_2$.

$$\mathcal{F} = \left\{ \max(f_1, f_2) : f_1 \in \mathcal{F}_1, f_2 \in \mathcal{F}_2 \right\}.$$

Given (x_1, x_2, \dots, x_n) .

$$|\mathcal{F}(x_1, x_2, \dots, x_n)| \leq |\mathcal{F}_1(x_1, \dots, x_n)| \cdot |\mathcal{F}_2(x_1, \dots, x_n)|.$$

$$\leq \left(\frac{en}{D_1}\right)^{D_1} \cdot \left(\frac{en}{D_2}\right)^{D_2}$$

$$\leq (en)^{D_1+D_2}$$

$$D = VC(\mathcal{F})$$

$$2^D \leq (e \cdot D)^{D_1+D_2}$$

$$D \leq C \cdot (D_1 + D_2).$$

eg. Linear threshold function.

$$\mathcal{F} = \left\{ \mathbb{1}_{\{\theta^T x \geq 0\}} : \theta \in \mathbb{R}^d \right\}$$

$(x_1, x_2, \dots, x_{d+1})$ \forall binary sequence b_1, b_2, \dots, b_{d+1}

$$\exists \theta, \text{ s.t. } b_i = \mathbb{I}\{\theta^T x_i \geq 0\} \quad (\forall i)$$

On the other hand $\exists \lambda_1, \lambda_2, \dots, \lambda_{d+1}$ not all 0.

$$\sum_{i=1}^{d+1} \lambda_i x_i = 0.$$

$$b_i = \begin{cases} 1 & \lambda_i > 0 \\ 0 & \text{otherwise.} \end{cases}$$

$$\lambda_i \theta^T x_i \geq 0. \quad \forall i$$

$$\sum_{i=1}^{d+1} \lambda_i \theta^T x_i = 0$$

$$\text{So } \lambda_i \theta^T x_i = 0 \quad (\forall i).$$

Naïve application of VC dim

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \right]$$

$$\leq C \sqrt{\frac{\log |\mathcal{F}(x_1, \dots, x_n)|}{n}}$$

$$\leq C \sqrt{\frac{D \log(n/D)}{n}}$$

suboptimal

Thm (VC dim bound for covering/packing). (\mathcal{F} binary)

$$\forall \epsilon, \log N(\epsilon; \mathcal{F}, \|\cdot\|_{L^2(Q)}) \leq C \cdot VC(\mathcal{F}) \cdot \log\left(\frac{1}{\epsilon}\right).$$

Proof — f_1, f_2, \dots, f_N $\max \epsilon$ -packing of \mathcal{F} under $L^2(Q)$.

$$\epsilon^2 \leq \|f_i - f_j\|_{L^2(Q)}^2 = \int_{\mathcal{X}} (f_i(x) - f_j(x))^2 dQ(x)$$

$$= Q(f_i(x) \neq f_j(x)) \quad (\forall i, j \in [N]).$$

Draw i.i.d. samples $x_1, x_2, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} Q$.

$$A = \left\{ (f_i(x_1), f_i(x_2), \dots, f_i(x_n)) : i \in [N] \right\}$$

$\forall i, j \in [N]$

$$P\left((f_i(x_1), f_i(x_2), \dots, f_i(x_n)) = (f_j(x_1), f_j(x_2), \dots, f_j(x_n)) \right)$$

$$\leq (1 - \epsilon^2)^n \leq \exp(-c\epsilon^2 n)$$

$\cdot P\left(\exists i, j \in [N] \text{ s.t. } f_i, f_j \text{ not distinguishable under } (x_1, \dots, x_n) \right)$

$$\leq \binom{N}{2} \cdot \exp(-c\epsilon^2 n) \leq N^2 \exp(-c\epsilon^2 n).$$

$$n = \frac{4 \log N}{\epsilon^2} \quad P(\dots) \leq \frac{1}{2}.$$

$\exists x_1, x_2, \dots, x_n$ s.t. $(f_i)_{i \in [N]}$ are pairwise distinguishable on (x_1, \dots, x_n) .

$$|A| = N \leq \left(\frac{en}{D}\right)^D \quad (D = VC(F))$$

$$N \leq \left(\frac{4e \log N}{D \epsilon^2}\right)^D$$

$$\log N \leq D \cdot \left\{ c + \log \log N + \log \left(\frac{1}{\epsilon}\right) \right\}$$

Solving for N ,

$$\log N \leq c \cdot D \log \left(\frac{1}{\epsilon}\right).$$

\mathcal{F} = set of linear classifiers in \mathbb{R}^d .

$$\text{gen err} \leq \mathbb{E} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E} f(X) \right| \right]$$

$$\leq \frac{C}{\sqrt{n}} \int_0^1 \sqrt{\log \sup_Q N(\epsilon; \mathcal{F}, \|\cdot\|_{L_2(Q)})} d\epsilon.$$

$$\leq \frac{C}{\sqrt{n}} \int_0^1 \sqrt{c' \cdot d \log(1/\epsilon)} d\epsilon$$

$$= C'' \sqrt{\frac{d}{n}}$$

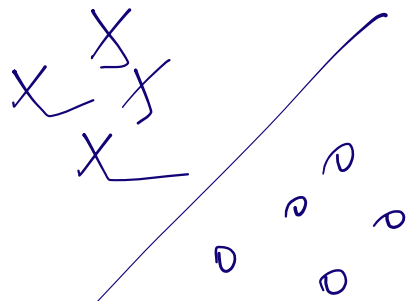
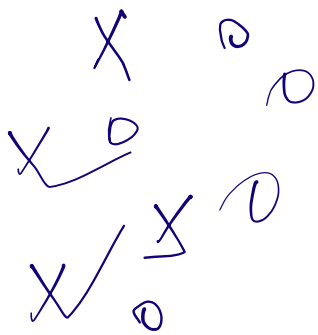
P. Bouchitté

S. Mendelson

65

"Local Rademacher Complexities"

(Tight when data are not separable)



Def: VC subgraph dimension

$$\mathcal{F} := \{x \mapsto \mathbb{1}_{\{t \leq f(x)\}} : f \in \mathcal{F}\}$$

$$VC(\mathcal{F}) := VC(\{ \mathbb{1}_{\{t \leq f(x)\}} : f \in \mathcal{F} \})$$



equivalently. largest D s.t. $\exists (x_i, t_i)_{i=1}^D$

$$\left\{ \mathbb{1}_{\{t_i \geq f(x_i)\}} \right\}_{i=1}^D \in \mathcal{F}$$

achieves all the binary configurations.

eg. $\mathcal{F} := \{x \mapsto \theta^T x : \theta \in \mathbb{R}^d\}$

$$t_i \geq f_\theta(x_i) \iff \begin{bmatrix} -t_i \\ x_i \end{bmatrix}^T \begin{bmatrix} 1 \\ \theta \end{bmatrix} \geq 0.$$

$$VC(\mathcal{F}) \leq d+1.$$

eg. φ is strictly increasing

$$\mathcal{F} := \{x \mapsto \varphi(\theta^T x) : \theta \in \mathbb{R}^d\}$$

$$t_i \geq f(x_i) \iff \begin{bmatrix} \varphi^{-1}(t_i) \\ x_i \end{bmatrix} \begin{bmatrix} 1 \\ -\varphi \end{bmatrix} \geq 0.$$

Thm. $\sup_Q N(\varepsilon \|F\|_{L^2(Q)}; \mathcal{F}, L^2(Q)) \leq \left(\frac{C}{\varepsilon}\right)^{C_{VC}(\mathcal{F})}$
 where F is envelope for \mathcal{F} . ($\forall \varepsilon > 0$)

Proof: $\|f - g\|_{L^2(Q)}^2 = \int_X |f(x) - g(x)|^2 dQ(x)$

$$\leq 2 \int F(x) |f(x) - g(x)| dQ(x)$$

$$\leq 2 \iint_{(x,t): |t| \leq F(x)} F(x) \cdot |1_{t \leq f(x)} - 1_{t \leq g(x)}| dQ(x) dt$$

$$\left(f(x) - g(x) = \int_{-F(x)}^{F(x)} (1_{t \leq f(x)} - 1_{t \leq g(x)}) dt \right)$$

(Cauchy Schwarz)

$$\leq 2 \cdot \sqrt{\iint_{|t| \leq F(x)} (\sqrt{F(x)})^2 dt dQ(x)}$$

$$\cdot \sqrt{\iint_{|t| \leq F(x)} (\sqrt{F(x)})^2 \cdot (1_{t \leq f(x)} - 1_{t \leq g(x)})^2 dt dQ(x)}$$

$$= 2\sqrt{2} \cdot \|F\|_{L^2(Q)}$$

$$\sqrt{\iint_{|t| \leq F(x)} (1_{t \leq f(x)} - 1_{t \leq g(x)})^2 \cdot F(x) dt dQ(x)}$$

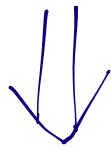
$d\tilde{Q}(x,t) \propto F(x) dt dQ(x)$
 on the domain $\{(t,x) : |t| \leq F(x)\}$

Normalizing const $= \iint_{|t| \leq F(x)} F(x) dt dQ(x) = 2 \cdot \|F\|_{L^2(Q)}^2$

So we have

$$\|f - g\|_{L^2(Q)}^2 \leq 4 \cdot \|F\|_{L^2(Q)}^2 \cdot \sqrt{\iint (1_{t \leq f(x)} - 1_{t \leq g(x)})^2 d\tilde{Q}(x,t)}$$

ϵ -covering of $\{(t,x) \rightarrow 1_{t \leq f(x)} : f \in \mathcal{F}\}$
 under \tilde{Q}

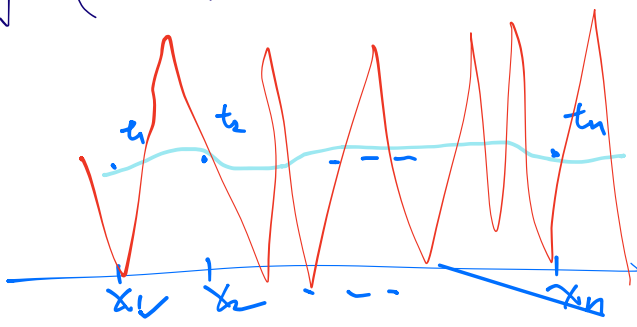


$2 \cdot \|F\|_{L^2(Q)} \cdot \sqrt{\epsilon}$ -covering of \mathcal{F} under Q .

Nonparametric classes usually have infinite

VC (subgraph) dimension.

Def: (ϵ -fat-shattering dimension)



$\text{fat}_\varepsilon(\mathcal{F})$ is largest D s.t.
 $\exists (x_i; t_i)_{i=1}^D$ ε -shattered by \mathcal{F} .

i.e.
 \forall binary seq b_1, b_2, \dots, b_D
 $\exists f \in \mathcal{F}$

$f(x_i) \geq t_i + \varepsilon$ when $b_i = 1$
 $f(x_i) < t_i$ when $b_i = 0$.

examples - future lectures on nonpara

Thm (Mendelson - Vershynin). \mathcal{F} bdd by 1

$$\sup_Q M(\varepsilon; \mathcal{F}, L^2(Q)) \leq \left(\frac{1}{\varepsilon}\right) c_1 \text{fat}_{c_2 \varepsilon}(\mathcal{F})$$

(where c_1, c_2 are universal consts)

Thm (Radelson - Vershynin) (Ann. Math.)
 under some mild assumptions

$$\sup_Q M(\varepsilon; \mathcal{F}, L^2(Q)) \leq \exp(c_1 \text{fat}_{c_2 \varepsilon}(\mathcal{F}))$$

Original problem

$$\hat{\theta}_n = \arg \min \frac{1}{n} \sum_{i=1}^n \ell(\theta; x_i) =: \hat{L}_n(\theta)$$

$$\theta^* = \arg \min \mathbb{E}[\ell(\theta; X)] =: L(\theta)$$

$$L(\hat{\theta}_n) - L(\theta^*) \leq 2 \cdot \sup_{\theta \in \Theta} |\hat{L}_n(\theta) - L(\theta)| \lesssim \sqrt{\frac{VC}{n}}$$

Too conservative. (at least locally around θ^*)

$$\|\hat{\theta}_n - \theta^*\| \lesssim \left(\frac{VC}{n}\right)^{1/4}$$

Thm (main convergence rate thm)

Assume $L(\theta) - L(\theta^*) \geq \|\theta - \theta^*\|^2$.

Suppose that

$$\mathbb{E} \left[\sup_{\substack{\theta \in \Theta \\ \|\theta - \theta^*\| \leq u}} (P_n - P)(\ell_\theta - \ell_{\theta^*}) \right] \leq \phi_n(u).$$

satisfying $\phi_n(cx) \leq C^\alpha \phi_n(x)$

(for any $c > 1, x > 0$) for some $\alpha < 2$.

Then, for any δ_n that satisfies $\phi_n(\delta_n) \leq \delta_n^2$

$\forall \epsilon > 0, \exists C_\epsilon > 0$, st. $\|\hat{\theta}_n - \theta^*\| \leq C_\epsilon \cdot \delta_n$ w.p. $1 - \epsilon$.

Proof. $\|\hat{\theta}_n - \theta^*\|^2 \leq L(\hat{\theta}_n) - L(\theta^*) \leq (P_n - P)(\ell_{\theta^*} - \ell_{\hat{\theta}_n})$.

$$\mathbb{P}(\|\hat{\theta}_n - \theta^*\| \geq 2^M \delta_n)$$

$$= \sum_{j>M} \mathbb{P}(2^{j-2} \delta_n \leq \|\hat{\theta}_n - \theta^*\| < 2^j \delta_n)$$

Each term $\leq \mathbb{P}(\|\hat{\theta}_n - \theta^*\| \leq 2^j \delta_n, |(P_n - P)(l_{\theta^*} - l_{\hat{\theta}_n})| \geq 2^{j-2} \delta_n^2)$

$$\leq \mathbb{P}\left(\sup_{\theta \in \Theta, \|\theta - \theta^*\| \leq 2^j \delta_n} |(P_n - P)(l_{\theta^*} - l_{\theta})| \geq 2^{j-2} \delta_n^2\right)$$

$$\leq \frac{1}{2^{j-2} \delta_n^2} \mathbb{E} \left[\sup \dots \right]$$

$$\leq \frac{1}{2^{j-2} \delta_n^2} \phi_n(2^j \delta_n)$$

$$\leq \frac{4 - \phi_n(\delta_n)}{\delta_n^2} \cdot 2^{(\alpha-2)j}$$

Summing up.

$$\mathbb{P}(\|\hat{\theta}_n - \theta^*\| \geq 2^M \delta_n)$$

$$\leq \frac{4 \phi_n(\delta_n)}{\delta_n^2} \cdot \frac{2^{-(\alpha-2)M}}{1 - 2^{\alpha-2}} \rightarrow 0.$$

≤ 4