

VC-dim is useful for binary functions
 $(X \rightarrow \{0, 1\})$.

How about real-valued functions.

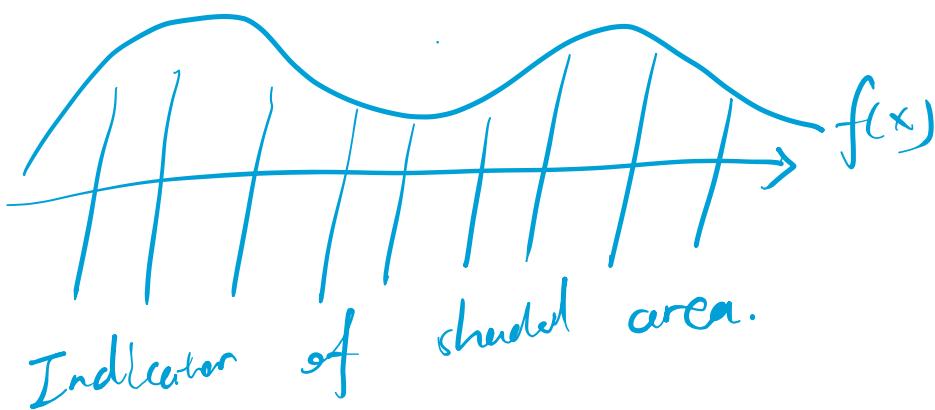
Def. VC subgraph dimension.

$$\mathcal{F} \subseteq (X \rightarrow \mathbb{R})$$

Convert to binary class

$$f(t, x) \mapsto \{t \leq f(x)\}$$

is a binary class



$\text{VC}(\mathcal{F}) := \text{VC-dim of the augmented subgraph function class}$

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

$$t \geq f_\theta(x) = \theta^T x \iff \begin{bmatrix} t \\ x \end{bmatrix} \cdot \begin{bmatrix} -1 \\ \theta \end{bmatrix} \leq 0.$$

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

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

for φ cs and strictly increasing φ .

$$t \geq f_\theta(x) \iff \begin{bmatrix} \varphi'(t) \\ x \end{bmatrix} \begin{bmatrix} 1 \\ -\theta \end{bmatrix} \geq 0.$$

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

Thm. $\sup_Q N(\varepsilon \|F\|_{L^2(Q)}; \mathcal{F}, L^2(Q)) \leq \left(\frac{C_1}{\varepsilon}\right)^{C_2 \cdot VC(\mathcal{F})}$

where C_1, C_2 are universal constants.

F is an envelop function of class \mathcal{F} .

Remark: this is exactly the covering/packing #
· there appeared in main empirical process bound.

Recall:

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_i^n f(X_i) - \mathbb{E} f(X) \right| \right]$$

$$\leq C \sqrt{\frac{\mathbb{E}[F(X)^2]}{n}} \int_Q \sqrt{\log \sup_Q N(\varepsilon \|F\|_{L^2(Q)}; \mathcal{F}, L^2(Q))} d\varepsilon$$

$$\leq C \cdot \sqrt{\frac{E(F(x)^2)}{n}} \cdot \int_0^1 \sqrt{C_{\epsilon} VCF} \cdot \log\left(\frac{C}{\epsilon}\right) d\epsilon$$

$$\leq C' \cdot \|F\|_{L^2(P)} \cdot \sqrt{\frac{VC(F)}{n}}.$$

Proof of "VC subgraph \rightarrow covering/packing":

Idea: convert covering/packing of subgraph class
to that of F itself.

For $f, g \in F$

$$\begin{aligned} \|f - g\|_{L^2(Q)}^2 &= \int_X |f(x) - g(x)|^2 dQ(x) \\ &\leq 2 \int_X F(x) \cdot |f(x) - g(x)| dQ(x). \end{aligned}$$

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

So we have

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

$$\begin{aligned}
&\leq 2 \cdot \underbrace{\iint_{|t| \leq F(x)} (\sqrt{F(x)})^2 dt dQ(x)}_{\cdot} \\
&\quad \cdot \underbrace{\iint_{|t| \leq F(x)} (\sqrt{F(x)})^2 (1_{|t| \leq f(x)} - 1_{|t| \leq g(x)})^2 dt dQ(x)}_{\cdot} \\
&= 2\sqrt{2} \cdot \|F\|_{L^2(Q)} \cdot \iint_{|t| \leq F(x)} (1_{|t| \leq f(x)} - 1_{|t| \leq g(x)})^2 \cdot F(x) dt dQ(x)
\end{aligned}$$

Define $\hat{dQ}(x, t) \propto F(x) dt dQ(x)$.
on the domain $\{t, x : |t| \leq F(x)\}$.

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

$$\text{So } \hat{dQ}(x, t) = \frac{F(x) dt dQ(x)}{2 \|F\|_{L^2(Q)}^2}$$

Putting them together,

$$\|f - g\|_{L^2(Q)}^2 \leq 4 \cdot \|F\|_{L^2(Q)}^2 \cdot \iint_{|t| \leq F(x)} (1_{|t| \leq f(x)} - 1_{|t| \leq g(x)})^2 d\hat{Q}$$

So we have

$$N\left(2\|F\|_{L^2(Q)}\sqrt{\varepsilon}; F, L^2(Q)\right) \leq N\left(\varepsilon; \{I_{\{t \leq f(x)\}} : f \in F\}, L^2(\tilde{Q})\right).$$

$$\leq \left(\frac{C}{\varepsilon}\right)^{VC(F)}$$

You can take an $L^2(\tilde{Q})$ covering of

$\{I_{\{t \leq f(x)\}} : f \in F\}$, and map it

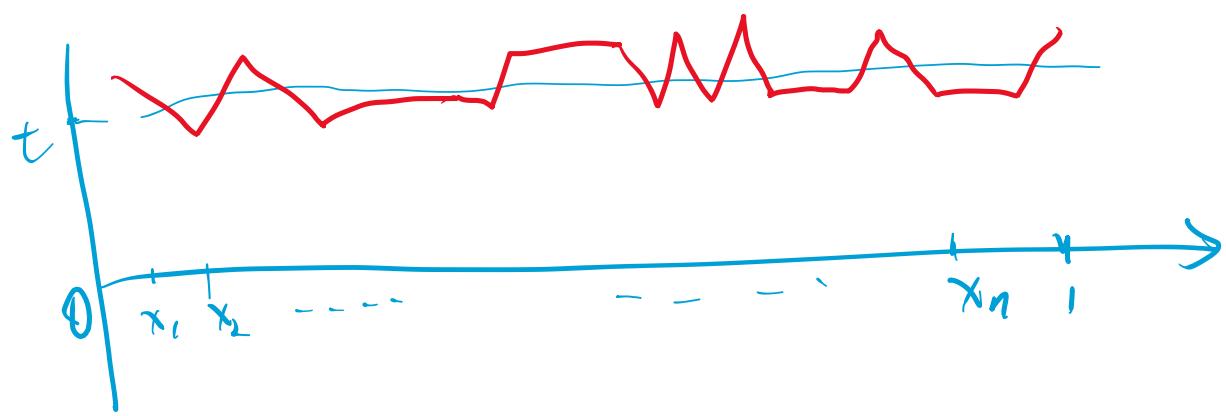
to $L^2(Q)$ covering of F .

How about infinite VC dim?

Motivating e.g.

$$F = \left\{ f: [0,1] \rightarrow [0,1] \mid |f(x) - f(y)| \leq |x-y|, \forall x, y \right\}$$

Any possible label sequence can be shattered.



$t_1 = t_2 = \dots = t_n = t$.
Any binary seq can be realised

$$VC(\mathcal{F}) = +\infty.$$

Solution: fat-shattering dim.

Def (ε -fat-shattering dim)

$\text{fat}_\varepsilon(\mathcal{F})$ is largest D s.t.

$\exists (x_i, t_i)_{i=1}^D$ that is ε -shattered by \mathcal{F} .

" ε -shattered": A binary seq b_1, b_2, \dots, b_D

$\exists f \in \mathcal{F}$. s.t. $f(x_i) \begin{cases} < t_i, & b_i = 0 \\ \geq t_i + \varepsilon, & b_i = 1 \end{cases}$

"Scale-sensitive combinatorial dim".

Why useful:

Thm (Mendelson- Vershynin)

If \mathcal{F} is unif bdd by 1,

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

So that by Dudley chaining integral, we get

$$\sup_{f \in \mathcal{F}} |(P_n - P)f| \leq \frac{c}{\sqrt{n}} \int_0^1 \sqrt{\text{fat}_{\varepsilon}(\mathcal{F}) \cdot \log(\frac{1}{\varepsilon})} d\varepsilon.$$

(Improved version by Rudelson- Vershynin)

under some mild assumptions

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

so that we have

$$\sup_{f \in \mathcal{F}} |(P_n - P)f| \leq \frac{c}{\sqrt{n}} \int_0^1 \sqrt{\text{fat}_{\varepsilon}(\mathcal{F})} d\varepsilon.$$

Fat-shattering dim of Lip func (in 1D)

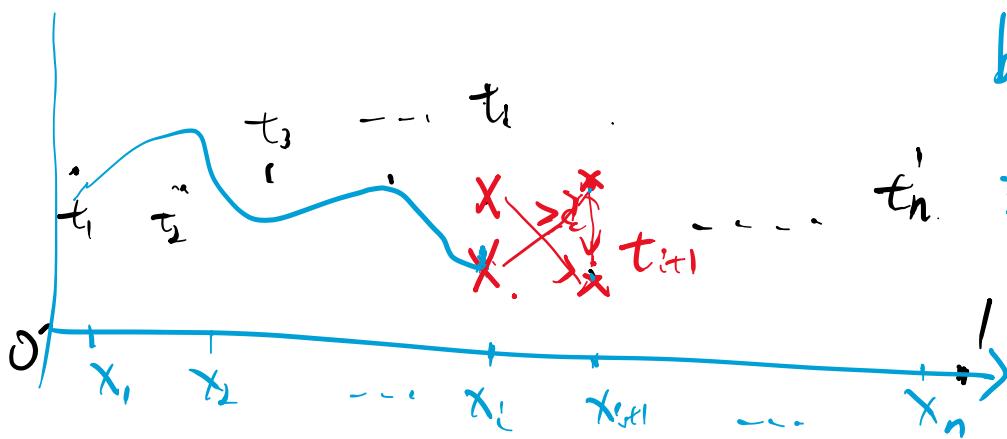
Suppose $(x_i, t_i)_{i=1}^n$ being ε -shattered.

(Assume $x_1 < x_2 < \dots < x_n$ w.l.o.g.).

A binary seq

b_1, b_2, \dots, b_n

$\exists f_b \in \mathcal{F}$.



s.t.

$$\begin{cases} f_b(x_i) < t_i & (b_i = 0) \\ f_b(x_{i+1}) \geq t_{i+1} + \varepsilon & (b_{i+1} = 1) \end{cases}$$

Fix b_1, b_2, \dots, b_{i-1}

Let f_+ be the function obtained by $b_i = 0, b_{i+1} = 1$

f_- --- $b_i = 1, b_{i+1} = 0$.

$$f_+(x_i) \geq t_i + \varepsilon, \quad f_+(x_{i+1}) \leq t_{i+1}$$

$$f_-(x_i) \leq t_i, \quad f_-(x_{i+1}) \geq t_{i+1} + \varepsilon$$

By Lip condition, $|f_t(x_i) - f_t(x_{i+1})| \leq |x_{i+1} - x_i|$

$$|f_-(x_i) - f_-(x_{i+1})| \leq |x_{i+1} - x_i|.$$

so we have

$$\begin{aligned}|x_{i+1} - x_i| &\geq \max\{t_i + \varepsilon - t_{i+1}, t_{i+1} + \varepsilon - t_i\} \\&= \varepsilon + |t_i - t_{i+1}| \geq \varepsilon.\end{aligned}$$

so we have $n \leq \frac{1}{\varepsilon}$.

and therefore $\text{fat}_\varepsilon(F) \leq \frac{1}{\varepsilon}$.

Original problem: $L_n(\theta) = \frac{1}{n} \sum_{i=1}^n l_\theta(x_i)$.

$$\hat{\theta}_n := \arg \min L_n(\theta) \quad \theta^* = \arg \min L(\theta).$$

$$L(\hat{\theta}_n) - L(\theta^*) \leq 2 \sup_{\theta \in \mathcal{H}} |L_n(\theta) - L(\theta)|.$$

(e.g. $\leq C \sqrt{\frac{VC(F)}{n}}$)

sup over \mathcal{H} can be too conservative.

$\hat{\theta}_n$ may be close to θ^*

If we know this a priori,

then we can probably take

\sup over $\{\theta \in \Theta \text{ neighbors of } \theta^*\}$.

Chicken-egg problem?

Thm (Localization).

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

Suppose that we can show

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

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

($A \ C > 1, x > 0$) for some $\alpha < 2$.

Then for smaller δ_n satisfying $\phi_n(\delta_n) \leq f_n^2$

$\forall \varepsilon > 0, \exists$ constant $C_\varepsilon > 0$ (depending only on ε)

s.t. $\|\hat{\theta}_n - \theta^*\| \leq C_\varepsilon \cdot \delta_n$ w.p. $1 - \varepsilon$.

δ_n = "critical radius".

Proof: $L(\hat{\theta}_n) - L(\theta^*) \geq \|\hat{\theta}_n - \theta^*\|^2.$

$$\begin{aligned}
&= L(\hat{\theta}_n) - L(\hat{\theta}_n) + \underbrace{L_n(\hat{\theta}_n) - L_n(\theta^*)}_{\leq 0} + L(\theta^*) - L(\theta^*) \\
&\leq (L_n(\theta^*) - L_n(\hat{\theta}_n)) - (L(\theta^*) - L(\hat{\theta}_n)) \\
&= (P_n - P)(\ell_{\theta^*} - \ell_{\hat{\theta}_n}).
\end{aligned}$$

Consider the tall prob

$$\begin{aligned}
&\mathbb{P}(\|\hat{\theta}_n - \theta^*\| \geq 2^M f_n) \\
&= \sum_{j > M} \mathbb{P}(2^{j-1} f_n \leq \|\hat{\theta}_n - \theta^*\| \leq 2^j f_n).
\end{aligned}$$

$$\begin{aligned}
\text{Each term} &\leq \mathbb{P}(\|\hat{\theta}_n - \theta^*\| \leq 2^j f_n, |(P_n - P)(\ell_{\theta^*} - \ell_{\hat{\theta}_n})| \geq 2^{j+2} f_n^2) \\
&\leq \mathbb{P}\left(\sup_{\theta \in \Theta} |(P_n - P)(\ell_{\theta^*} - \ell_{\theta})| \geq 2^{j+2} f_n^2 \mid \|\theta - \theta^*\| \leq 2^j f_n\right) \\
&\leq \frac{1}{2^{2j+2} f_n^2} \cdot \mathbb{P}_n(2^j f_n) \\
&\leq \frac{4 \mathbb{P}_n(f_n)}{f_n^2} \cdot 2^{(\alpha-2)j}.
\end{aligned}$$

Sum them up, we have

$$\begin{aligned} \mathbb{P}(\|\hat{\theta}_n - \theta^*\| \geq 2^M f_n) &\leq \frac{4\phi(f_n)}{f_n^2} \cdot \frac{2^{-(\alpha-2)M}}{1 - 2^{\alpha-2}} \\ &\leq 4 \end{aligned}$$

choose M
large enough,
make it $\leq \frac{\varepsilon}{4}$.

Remark: bad dependence on tail prob (ε)

usually sharp in terms of n
and complexity of f_n .

e.g. "Regular" M-estimators $\mathbb{D} = \mathbb{R}^d$.

$$|f_{\theta_1}(x) - f_{\theta_2}(x)| \leq M(x) \cdot \|\theta_1 - \theta_2\|_2$$

$$\text{Assume that } L(\theta) - L(\theta^*) \geq \|\theta - \theta^*\|_2^2$$

then we can compute

$$\mathbb{E} \left[\sup_{\substack{\theta \in \mathbb{R}^d \\ \|\theta - \theta^*\|_2 \leq u}} |(P_n - P)(f_\theta - f_{\theta^*})| \right].$$

$$\left(\forall \theta \text{ s.t. } \|\theta - \theta^*\|_2 \leq u, \quad |f_\theta(x) - f_{\theta^*}(x)| \leq u \cdot M(x) \right)$$

$$\leq C \frac{u}{\sqrt{n}} \|M\|_{L^2(P)}.$$

$$= \int_0^1 \sqrt{\log \sup_Q N(u \cdot \delta \cdot \|M\|_{L^2(Q)}; F_u, \| \cdot \|_{L^2(Q)})} d\delta.$$

$$(F_u := \{\theta - \theta^* : \|\theta - \theta^*\|_b \leq u\})$$

$$N(u \cdot \delta \cdot \|M\|_{L^2(Q)}; F_u, \| \cdot \|_{L^2(Q)})$$

$$(\|\theta_i - \theta_j\|_b \leq \varepsilon, \Rightarrow \|\ell_{\theta_i} - \ell_{\theta_j}\|_{L^2(Q)} \leq \varepsilon \|M\|_{L^2(Q)})$$

$$\leq N(u \delta; \{\theta : \|\theta - \theta^*\|_b \leq u\}, \| \cdot \|_b)$$

$$\leq \left(\frac{C}{\delta}\right)^d.$$

and we can conclude

$$\int_0^1 \sqrt{\log \sup - \dots} d\delta \leq C' \sqrt{d}.$$

$$\phi_n(u) = C \sqrt{\frac{d}{n}} \cdot \|M\|_{L^2(P)} \cdot u.$$

$$\text{Solving fixed pt} \quad f_n = C \sqrt{\frac{d}{n}} \cdot \|M\|_{L^2(P)}.$$