

Recap : covering / packing #.

$$\text{Thm (duality)} \quad M(K; p, 2\epsilon) \stackrel{\text{covering}}{\leq} N(K; p, \epsilon) \stackrel{\text{(i)}}{\leq} M(K; p, \epsilon) \stackrel{\text{packing}}{\leq}$$

↑
set ↑ metric ↑ radius

Proof: i) Let $\{\theta_j\}_{j=1}^M$ be 2ϵ -packing of K

$\{\psi_j\}_{j=1}^N$ be ϵ -covering of K .
 $\text{(want to show } M \leq N\text{).}$

$\forall \theta_i, (i=1, 2, \dots, M), \exists j \text{ s.t. } \theta_i \in B_p(\psi_j, \epsilon)$.

(more than one - choose arbitrarily).

Suppose if $\exists \theta_1, \theta_2 \in \{\theta_i\}$

s.t. $\theta_1, \theta_2 \in B_p(\psi_j, \epsilon)$

then $p(\theta_1, \theta_2) \leq p(\theta_1, \psi_j) + p(\psi_j, \theta_2) \leq 2\epsilon$

Contradiction ($\{\theta_i\}_{i \in [M]}$ is 2ϵ -packing).

So different θ 's map to different ψ .

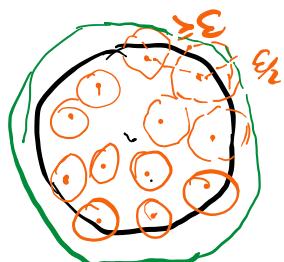
$$M \leq N,$$

(!!). Let $\{\theta_j\}_{j=1}^M$ be a maximal ε -packing.
 By maximality, $\forall \theta \in K, \exists j \in [M]$ s.t. $p(\theta, \theta_j) \leq \varepsilon$
 (otherwise, θ can be added to the packing).
 So $\{\theta_j\}_{j=1}^M$ is also an ε -covering.

$$\underbrace{N(K; p, \varepsilon)}_{\text{minim. } \varepsilon\text{-covering}} \leq \underbrace{M}_{\text{an } \varepsilon\text{-covering}}$$

e.g. (Euclidean space).
 $K = \{\theta \in \mathbb{R}^d : \|\theta\|_2 \leq 1\}, p = \|\cdot\|_2$.

• Upper bound:



$$\underbrace{\bigcup_{i=1}^M B(\theta_i, \varepsilon/2)}_{\text{disjoint.}} \subseteq B(0, 1 + \frac{\varepsilon}{2})$$

$$\sum_{j=1}^M \text{Vol}(B(\theta_j, \varepsilon/2)) \leq \text{Vol}\left(B(0, 1 + \frac{\varepsilon}{2})\right)$$

$$M \leq \left(1 + \frac{2}{\varepsilon}\right)^d.$$

• Lower bound: suppose

$$K \subseteq \bigcup_{j=1}^N B(\psi_j, \varepsilon).$$

$$\text{Vol}(K) \leq \sum_{j=1}^N \text{Vol}(B(\psi_j, \varepsilon)). \quad \text{So } N \geq \left(\frac{1}{\varepsilon}\right)^d.$$

By duality thm:

$$\left(\frac{1}{\varepsilon}\right)^d \leq N(K; \| \cdot \|_2, \varepsilon) \leq \left(1 + \frac{2}{\varepsilon}\right)^d.$$

Back to statistics.

Consider M-estimation

$$\hat{\theta}_n = \underset{\theta \in \Theta}{\operatorname{arg\,min}} F_n(\theta)$$

$$\text{where } F_n(\theta) := \frac{1}{n} \sum_{i=1}^n f(\theta; x_i).$$

from last lecture, we know $\hat{\theta}_n \xrightarrow{P} \theta^*$
is implied by

- θ^* is unique minimizer of $F(\theta) := \mathbb{E}[f(\theta; X)]$
- (F is cts) — (implies ε - δ condition in last lecture).
- $\sup_{\theta \in \Theta} |F_n(\theta) - F(\theta)| \xrightarrow{P} 0$. (ULLN).

Thm (Wald) Suppose Θ is a compact set, and

$$(i) \mathbb{E}\left[\sup_{\theta \in \Theta} |f(\theta; X)|\right] < \infty. \quad (\text{DCT-type condition})$$

(ii) $f(\cdot; X)$ is cts (w.r.t. first argument).

then ULLN $\sup_{\theta \in \Theta} |F_n(\theta) - F(\theta)| \xrightarrow{P} 0$ holds

(Combining w/ uniqueness of θ^* , it leads to consistency).

Proof. since $f(\cdot; \theta)$ is v.s., F iscts on Θ

so F is uniformly v.s.

Fix $\varepsilon > 0$, $\exists \delta$ s.t. $\forall \theta, \theta' \in \Theta, \rho(\theta, \theta') \leq \delta$,
we have $|F(\theta) - F(\theta')| \leq \varepsilon$.

(Recall: we're able to prove ULLN for finite Θ ,

Compact Θ : use discretization)

Let $\{\theta_j\}_{j=1}^N$ be a δ -covering of Θ .

$$\begin{aligned} & \sup_{\theta \in \Theta} |F_n(\theta) - F(\theta)| \\ & \leq \max_{i \in [N]} \sup_{\substack{\theta \in \Theta \\ \rho(\theta, \theta_i) \leq \delta}} \left(|F_n(\theta) - F_n(\theta_i)| + |F(\theta) - F(\theta_i)| \right) \\ & \quad \text{need technicality.} \\ & \quad \text{①} \quad \text{②} \\ & \quad \text{depends on } i \text{ but not } \theta. \end{aligned}$$

max_{i ∈ [N]} |F_n(θ_i) - F(θ_i)| →_P 0
(for fixed covering).

First term

$$\begin{aligned} & = \max_{i \in [N]} \sup_{\theta \in B_p(\theta_i, \delta)} \left| \frac{1}{n} \sum_{j=1}^n (f(\theta; x_j) - f(\theta_i; x_j)) \right| \\ & \leq \max_{i \in [N]} \frac{1}{n} \sum_{j=1}^n \sup_{\theta \in B_p(\theta_i, \delta)} |f(\theta; x_j) - f(\theta_i; x_j)| \end{aligned}$$

$$\begin{aligned}
 & \xrightarrow{\text{(by LLN)}} \max_{\theta \in [N]} \mathbb{E} \left[\sup_{\theta \in B_p(\theta_j; \delta)} |f(\theta; x) - f(\theta_j; x)| \right] \\
 & \leq \sup_{\theta \in K} \mathbb{E} \left[\sup_{\rho(\theta'; \theta) \leq \delta} |f(\theta; x) - f(\theta'; x)| \right]
 \end{aligned}$$

Only need to show

$$\sup_{\theta \in K} \mathbb{E} \left[\sup_{\rho(\theta'; \theta) \leq \delta} |f(\theta; x) - f(\theta'; x)| \right] \xrightarrow{P} 0$$

(as $\delta \rightarrow 0$).
→ 0 ptwise.

$$M_\delta(\theta) := \sup_{\rho(\theta'; \theta) \leq \delta} |f(\theta; x) - f(\theta'; x)| \rightarrow 0.$$

$$|M_\delta(\theta)| \leq 2 \sup_{\theta \in \Theta} |f(\theta; x)| \in L^1.$$

So by DCT, $\mathbb{E}[M_\delta(\theta)] \rightarrow 0$ (by D).

By Dini's thm, convergence is uniform in $\theta \in \Theta$.

Remark:

- Only need weak conditions, works for general metric spaces
- metric is used in $\begin{cases} \text{compactness of } \Theta \\ \text{cts func } f(\cdot; x) \end{cases}$
- Quantitative bound?

For ②:

$$\begin{aligned} & \mathbb{P}\left(\max_{i \in [N]} |F_n(\theta_i) - F(\theta_i)| \geq \varepsilon\right) \\ & \leq \sum_{i \in [N]} \mathbb{P}(|F_n(\theta_i) - F(\theta_i)| \geq \varepsilon) \\ & \leq N \cdot \sup_{\theta \in \Theta} \mathbb{P}(|F_n(\theta) - F(\theta)| \geq \varepsilon) \end{aligned}$$

- If f is bdd, $\mathbb{P}(|\dots| \geq \varepsilon) \leq \exp(-cn\varepsilon^2)$,
in order to make $\mathbb{P}(\max |\dots| \geq \varepsilon) \leq \delta$.

we choose $\varepsilon \asymp \sqrt{\frac{\log(N/\delta)}{n}}$.

- If f is heavy-tailed. (e.g. $\mathbb{E}[f(\theta; X)^2] < \infty$)
naïve application of Chebychev ineq gives

$$N \cdot \mathbb{P}(|\dots| \geq \varepsilon) \leq O\left(\frac{N}{n\varepsilon^2}\right).$$

$$\varepsilon \asymp \sqrt{\frac{N}{n\delta}}.$$

This can be overcome by
"symmetrization method".

For ①: one-step discretization
is suboptimal.
use multi-level discretization.
"chaining"

"Symmetrization":

$$P_n f := \frac{1}{n} \sum_{i=1}^n f(x_i) \quad P_f := E[\bar{f}(X)].$$

want to bound $\sup_{f \in \mathcal{F}} |P_n f - P_f|$.

(we replace $f(0; X)$ with just $f(X)$).

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

where $(\varepsilon_i)_{i=1}^n$ are iid Rademacher r.v.'s indep of $(X_i)_{i=1}^n$

$$\left(\varepsilon_i = \begin{cases} 1 & \text{w.p. } \frac{1}{2} \\ -1 & \text{w.p. } \frac{1}{2} \end{cases} \right)$$

"Rademacher complexity" of \mathcal{F} .

Thm. $E \left[\sup_{f \in \mathcal{F}} |P_n f - P_f| \right] \leq 2 \cdot R_n(\mathcal{F}).$

Why useful? we can condition on data $(X_i)_{i=1}^n$

bound $\mathbb{E}\left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_i) \right| \middle| (X_i)_{i=1}^n\right]$.

easy to use concentration ineq.

Proof. $\mathbb{E}\left[\sup_{f \in \mathcal{F}} |P_n f - P f|\right]$

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

$$\stackrel{(Jensen)}{\leq} \mathbb{E}\left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(X_i) - f(X'_i)) \right| \right].$$

$$f(X_i) - f(X'_i) \stackrel{d}{=} \varepsilon_i (f(X_i) - f(X'_i))$$

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

$$\stackrel{(Jensen)}{\leq} 2 \cdot \mathbb{E}\left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_i) \right| \right].$$

For $A \subseteq \mathbb{R}^n$, want to bound

$$R_n(A) := \mathbb{E}\left[\sup_{a \in A} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i a_i \right| \right].$$

Complexity measure of set A .

"max correlation w/ a random direction".

related to Gaussian complexity
(replace Rademacher r.v. with $N(0, 1)$).

• Finite $A \subseteq \mathbb{R}^n$.

$\forall \lambda > 0$, by Jensen's ineq.

$$R_n(A) \leq \frac{1}{\lambda} \log \mathbb{E} \left[\exp \left(\max_{a \in A} \frac{\lambda}{n} \sum_{i=1}^n \varepsilon_i a_i \right) \right] \quad (\text{MGF})$$

$$\leq \frac{1}{\lambda} \log \left(\sum_{a \in A} \mathbb{E} \left[\exp \left(\frac{\lambda}{n} \sum_{i=1}^n \varepsilon_i a_i \right) \right] \right) \quad (\text{union bound}).$$

$$\leq \frac{1}{\lambda} \log \left(\sum_{a \in A} \exp \left(\frac{\lambda^2}{2n^2} \|a\|_2^2 \right) \right).$$

$$\leq \frac{1}{\lambda} \log |A| \cdot \exp \left(\frac{\lambda^2}{2n^2} \max_{a \in A} \|a\|_2^2 \right)$$

$$= \frac{\log |A|}{\lambda} + \frac{\lambda}{2n^2} \cdot \max_{a \in A} \|a\|_2^2.$$

(choose optimal λ)

$$= \max_{a \in A} \frac{\|a\|_2}{\sqrt{n}} \cdot \sqrt{\frac{2 \log |A|}{n}}.$$

(Alternatively, you may also use Hoeffding + union bound
to get the same result).

Discretization:

— One-step

$a^{(1)}, a^{(2)}, \dots, a^{(M)}$ δ -covering of A .

under $\|a\|_n := \sqrt{\frac{1}{n} \sum_{i=1}^n a_i^2}$.

$$\pi(a) := \underset{a^{(i)}}{\operatorname{argmin}} \|a - a^{(i)}\|_n.$$

$$\begin{aligned} R_n(A) &\leq \mathbb{E} \left[\sup_{a \in A} \left| \frac{1}{n} \epsilon^T (a - \pi(a)) \right| \right] + \mathbb{E} \left[\max_{j \in [M]} \frac{1}{n} \epsilon^T a^{(j)} \right] \\ &\quad \text{(Cauchy-Schwarz)} \\ &\leq \mathbb{E} \left[\sup_{a \in A} \|\epsilon\|_n \|a - \pi(a)\|_n \right] \leq \delta. \end{aligned}$$

$$\text{So } R_n(A) \leq \delta + \operatorname{diam}(A) \cdot \sqrt{\frac{2 \log M(A; \delta)}{n}}$$

e.g. $F = d$ -dim parametric model

then A is a d -dim manifold.

$$\text{choose } \delta = O\left(\frac{1}{n}\right). \quad M(A; \delta) \asymp \left(\frac{1}{\delta}\right)^d$$

$$\text{this leads to } R_n(A) \leq O\left(\sqrt{\frac{d \log n}{n}}\right)$$

Gap is even larger for nonparametric models.

Chaining: $A_0, A_1, A_2, \dots, A_m, \dots$

discrete approximations to A .

(approx err $\xrightarrow{m \rightarrow \infty} 0$)

$$\mathbb{E}\left[\sup \frac{1}{n} \sum \epsilon^T a\right] \leq \sum_{m=0}^{+\infty} \mathbb{E}\left[\sup_{a \in A} \left| \frac{1}{n} \sum (\pi_{m+1}(a) - \pi_m(a)) \right| \right]$$

where $\pi_m(a)$ is the best approx to a in A_m

$$a \in A : a = \sum_{m=0}^{+\infty} (\pi_{m+1}(a) - \pi_m(a))$$

$$(\pi_0(a) = 0)$$

Assuming that A_m is f_m -covering of A .

We can bound

$$\mathbb{E}\left[\sup_{a \in A} \left| \frac{1}{n} \sum (\pi_m(a) - \pi_{m+1}(a)) \right| \right]$$

$$\leq (f_m + f_{m+1}) \cdot \sqrt{\frac{2 \log |A_m| \cdot |A_{m+1}|}{n}} \| \pi_m(a) - \pi_{m+1}(a) \|_n \leq f_m + f_{m+1}$$

• $(\pi_m(a), \pi_{m+1}(a))$ takes value

in a finite set,

Cardinality $\leq |A_m| \cdot |A_{m+1}|$

Thm (Dudley's chaining).

\exists a universal const $c > 0$.

$$R_n(A) \leq \frac{c}{\sqrt{n}} \int_0^{+\infty} \sqrt{\log N(A; \|\cdot\|_n, \delta)} d\delta.$$

Proof: $D = \max_{a \in A} \|a\|_n$.

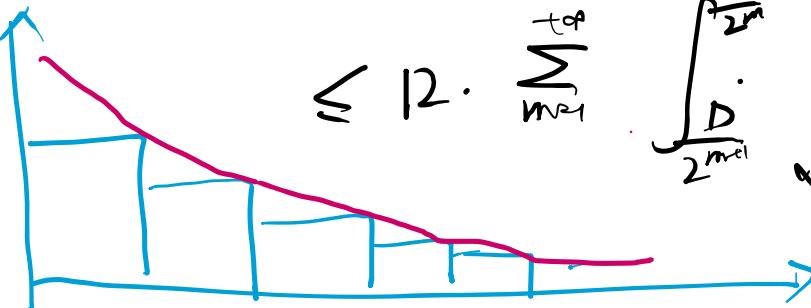
$A_0 = \{0\}$. $A_m :=$ minimal δ_m -covering

$$f_m = \frac{D}{2^m}.$$

$$(\delta_m + \delta_{m+1}) = \frac{3D}{2^m}$$

$$|A_m| \cdot |A_{m+1}| \leq |A_{m+1}|^2 = N(A; \|\cdot\|_n, \frac{D}{2^{m+1}})^2.$$

Substitute to the bound we derived.

$$\mathbb{E} \left[\sup_{a \in A} \left(\frac{1}{n} \epsilon^T a \right) \right] \leq \sum_{m=1}^{+\infty} \frac{3D}{2^m} \cdot \sqrt{\frac{4 \log N \left(\frac{D}{2^{m+1}} \right)}{n}} \\ \leq 12 \cdot \sum_{m=1}^{+\infty} \frac{D}{2^{m+1}} \cdot \sqrt{\frac{\log N(\delta)}{n}} d\delta.$$


Remarks:

- Dudley integral may diverge
of covering # upper bound.

In that case, we use different
upper bounds for diff f .

- "generic chaining" is tighter bounds
(known to be optimal),
but it's less convenient.

Back to function class.

Consider \mathcal{F} . assume \exists envelope function G .

$$|f(x)| \leq G(x) \quad (\forall f \in \mathcal{F})$$

(similar to Wald's consistency proof).

Thm. under above setup.

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} |P_n f - P f| \right]$$

$$\leq C \cdot \sqrt{\frac{\mathbb{E}[G(x)^2]}{n}} \cdot \int_0^1 \sqrt{\log \sup_Q N(\delta \cdot \|G\|_{L^2(Q)}; \mathcal{F}, L^2(Q))} d\delta.$$

"worst-case L^2 covering #":

Can be bounded in most cases
(see next lecture).