# Lecture 5: Finite-Sample Bounds for Minimum Distance Functionals

— Hw 1: extended to next Tuesday

## Recap

Concentration inequalities + union bound

$\Rightarrow$ bound supremum

• Maximum eigenvalue of random matrix ($\varepsilon$-net)
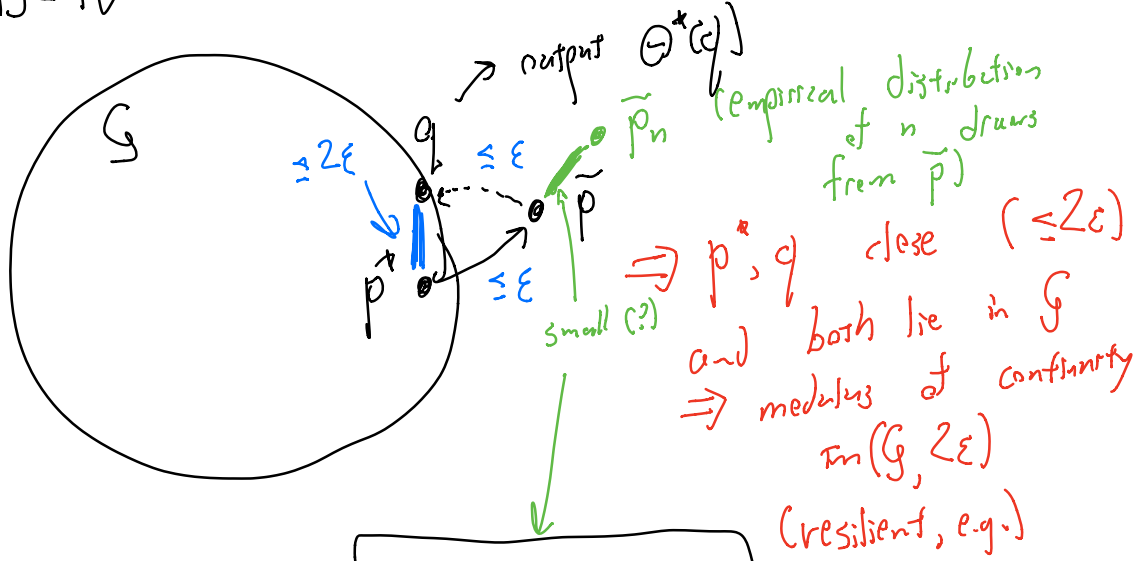
• VC dimension bound (symmetrization)

If $\mathcal{H}$ is a family of functions $f: \mathcal{X} \to \{0,1\}$, and $vc(\mathcal{H}) = d$, then w.p. $1-\delta$,

$$\sup_{f \in \mathcal{H}} \left| \nu_n(f) - \nu(f) \right| \leq \mathcal{O}\left( \sqrt{\frac{d + \log(1/\delta)}{n}} \right)$$

$\underbrace{\qquad}_{\widehat{\mathbb{P}}_n[f(x)=1]} \quad \underbrace{\qquad}_{\mathbb{P}(f(x)=1)}$

$\downarrow$
$\frac{1}{n} \sum_{i=1}^{n} f(x_i), \quad x_1, \ldots, x_n \sim \mathbb{P}$

Today: Finite-sample of MD functional

$D = TV$

output $\Theta^*(q)$

$G$

$\widetilde{p}_n$ (empirical distribution of $n$ draws from $\widetilde{p}$)

$q$ $\leq \varepsilon$

$\leq 2\varepsilon$

$\widetilde{p}$

$p^*$ $\leq \varepsilon$

small (?)

$\Rightarrow p^*, q$ close $(\leq 2\varepsilon)$

and both lie in $G$

$\Rightarrow$ modulus of continuity $m(G, 2\varepsilon)$ (resilient, e.g.)

hopefully? going to $0$ as $n \to \infty$

$TV(\widetilde{p}, \widetilde{p}_n) \rightrightarrows 1$ even as $n \to \infty$.

continuous

discrete

$TV = 1$ $\left(\begin{array}{l} \widetilde{p} = \text{Gaussian} \\ \widehat{p}_n: \end{array}\right.$

$\widetilde{p}_n$

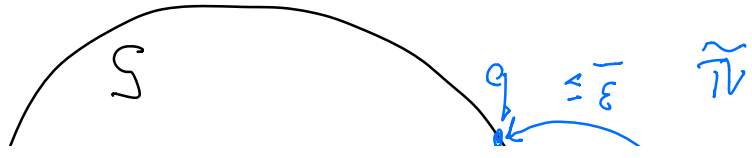$p^* \to \widetilde{p} \to \widetilde{p}_n$ | $p^* \to \widetilde{p}$

Issue. $TV(\bar{p}, \bar{p}_n)$ large $\Rightarrow$ can't directly apply $\Delta$-ineq.

Solutions. "relaxing the distance"

➀ Replace $TV$ w/ relaxed distance $\widetilde{TV}$
(today) where $\widetilde{TV}(\bar{p}, \bar{p}_n)$ actually is small

"expanding the set"

➁ (next time)



Lemma. Suppose $\widetilde{TV}$ is a distance s.t. $TV \geq \widetilde{TV}$ and
let $\bar{\varepsilon} = \varepsilon + TV(\bar{p}, \bar{p}_n)$. Then, the error of MD functional
for $\widetilde{TV}$ is at most $m(\mathcal{G}, \widetilde{TV}, 2\bar{\varepsilon})$.

$\leq 2\tilde{\varepsilon}$   $\leq \tilde{\varepsilon}$

$p^*$   $\tilde{p}$   $\tilde{p}_n$

$\Rightarrow p^*, y$  $2\tilde{\varepsilon}$ -close under $\widetilde{TV}$ distance

$\Rightarrow$ error bound by $m(g, \widetilde{TV}, 2\tilde{\varepsilon})$.

---

<u>Conclusion.</u>  Need: 

(a) $TV \geq \widetilde{TV}$

(1) $\widetilde{TV}(\tilde{p}, \tilde{p}_n) \to 0$

(2) need modulus to still be small

---

$\mathcal{H}$: family of fnc's $f: \mathcal{X} \to \mathbb{R}$

"variational representations"

$$\widetilde{TV}_{\mathcal{H}}(p, q) = \sup_{f \in \mathcal{H}, \tau \in \mathbb{R}} \left| \mathbb{P}_{X \sim p}[f(X) \geq \tau] - \mathbb{P}_{X \sim q}[f(X) \geq \tau] \right|$$

↑ bounded

$\sup |r(f) - r_n(f)|$

(a) $TV \geq \widetilde{TV}_{\mathcal{H}}$

$$TV(p, q) = \sup_E \left| \mathbb{P}_p[E] - \mathbb{P}_q[E] \right|$$

(1) $\widetilde{TV}_{\mathcal{H}}(\tilde{p}, \widehat{p}_n) \to 0$  as  $n \to \infty$  (or least intuitively)

formalize via VC-dim
argument

③ malalus → next

$$TV_1(p,q) = \sup_{\text{1-Lipschitz } f} |\mathbb{E}_p[f(x)] - \mathbb{E}_q[f(x)]|$$

Claim. For $(\beta, \xi)$-resilient distributions, malalus still bounded by $\beta$.

P •———$TV \le \xi$———• q

r •  $\exists r \le \frac{P}{1-\xi}, r \le \frac{q}{1-\xi}$

If $TV(p,q) \le \xi$, then

$\Rightarrow \mu(r) \approx \mu(p)$

$\mu(r) \approx \mu(q)$

$\triangle\text{-ineq.} \Rightarrow \mu(p) \approx \mu(q)$.

• P

• q

$\downarrow \qquad\qquad \downarrow$

$r_p \qquad\qquad\qquad r_q$

"mean of $p$ crosses the mean of $q$"

For any $f \in \mathcal{H}$, $\mathbb{E}_p[f(x)]$

Lemma. (Mean cross lemma) If $\widetilde{TV}_{\mathcal{H}}(p,q) \leq \varepsilon$, and

$f \in \mathcal{U}$, then $\exists \; r_p \leq \frac{p}{1-\varepsilon}, \; r_q \leq \frac{q}{1-\varepsilon}$, such that

$\mathbb{E}_{r_q}[f(x)] \leq \mathbb{E}_{r_p}[f(x)]$.

$\geq$

$r_p \leq \frac{p}{1-\varepsilon} \implies$ For all events $E$,

$r_p(E) \leq \frac{p(E)}{1-\varepsilon}$.

$\not\implies r_p(E) = \frac{p(E)}{1-\varepsilon} \implies r_p(\mathbb{R}^n) = \frac{1}{1-\varepsilon} \neq 1$

Pf of lemma, $\widetilde{TV}_{\mathcal{H}}(p,q) \leq \varepsilon$

$\implies \sup\limits_{\tau \in \mathbb{R}} \left| \mathbb{P}_p[f(x) \geq \tau] - \mathbb{P}_q[f(x) \geq \tau] \right| \leq \varepsilon$

$\widetilde{TV}_{\mathcal{H}}; \; \forall \; f \in \mathcal{H}$

CDF of $p, q$ $\varepsilon$-close



$p \qquad\qquad q$

$\varepsilon$

$\varepsilon \rightarrow 3rd[\text{?}]$

$p$

$f(x)$

$r_p$    $\Leftarrow r_q$    $r_q$

$r_p \Rightarrow$

For all $\tau$, $\mathbb{P}_{r_p}[f(x) \geq \tau] \geq \mathbb{P}_{r_q}[f(x) \geq \tau]$.

$\longrightarrow$ $r_p$ stochastically dominates $r_q$

$\Updownarrow$

$\longrightarrow \mathbb{E}_{r_q}[f(x)] \leq \mathbb{E}_{r_p}[f(x)]$

$$\mathbb{E}[Z] = \int_0^\infty \mathbb{P}[Z \geq \tau] \, d\tau$$

$\Leftarrow$                     stoch. dom.
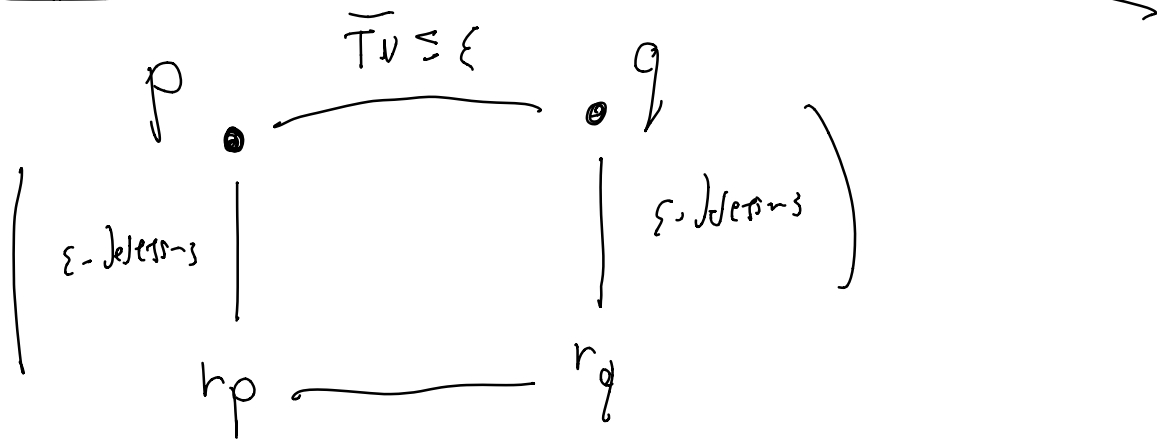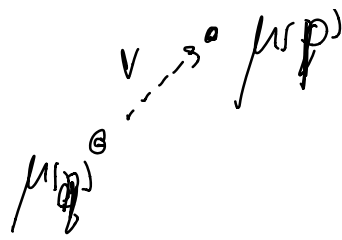
---

Back to bounding modulus.

Suppose $p, q$ are $(\beta, \varepsilon)$-resilient.

$\mathcal{H} = \{ x \mapsto \langle v, x \rangle \mid v \in \mathbb{R}^d \}$

If $\text{TV}_{\mathcal{H}}(p, q) \leq \varepsilon$.

$$v = \frac{\mu(p) - \mu(q)}{\|\mu(p) - \mu(q)\|_2}$$

$$\mu(q)^{\circ} \quad v \cdots \overset{\circ}{\longrightarrow} \mu(p)$$

$$\overline{Tv \leq \xi}$$

$$p \bullet \qquad \qquad \bullet q$$

$\varepsilon$-deletions $\Big|$ $\qquad$ $\Big|$ $\xi$-deletions

$$r_p \qquad \qquad r_q$$

$$\|\mu(p) - \mu(q)\|_2$$

$$= \langle v, \ \mu(p) - \mu(q) \rangle$$

$$= \langle \textcolor{blue}{\boxed{v}}, \mu(p) - \mu(r_p) \rangle \textcolor{blue}{\longrightarrow \leq \|\mu(p) - \mu(r_p)\|_2}$$

$$\textcolor{blue}{\leq \rho \ \text{(by resilience)}}$$

$$\textcolor{red}{+ \langle v, \mu(r_p) - \mu(r_q) \rangle} \textcolor{blue}{\leq 0}$$

$$\textcolor{blue}{= \mathbb{E}_{r_p}[\langle v, x \rangle] - \mathbb{E}_{r_q}[\langle v, x \rangle]}$$

$$\textcolor{blue}{\leq 0}$$

$$+ \langle v, \mu(r_q) - \mu(q) \rangle \textcolor{blue}{\longrightarrow \leq \rho}$$

$\leq 2\rho. \quad \Rightarrow \quad \| \mu(p) - \mu(q) \|_2 \leq 2\rho \quad \text{even for } \widetilde{TV}$

$\Rightarrow \quad \text{modulus of continuity bound.} \boxed{\text{fig}}$

---

**Thm.** If $p^t$ is $\left( \rho, 2\widetilde{\epsilon} \right)$-resilient

for $\widetilde{\epsilon} = \epsilon + \underbrace{\widetilde{TV}(\bar{p}, \widetilde{p}_n)}_{\lesssim \sqrt{\frac{d}{n}}}$, then

$\widetilde{TV}$ MD-functional recovers mean

w/ error $\leq 2\rho$, even in

finite samples.

$\boxed{\bigstar}$

---

$$\widetilde{TV}_{\mathcal{H}}(\bar{p}, \widetilde{p}_n) = \sup_{f \in \mathcal{H}, \tau \in \mathbb{R}} \left| \mathbb{P}_{\bar{p}}[f(x) \geq \tau] - \mathbb{P}_{\widetilde{p}_n}[f(x) \geq \tau] \right|$$

$$\mathcal{H}' = \left\{ \mathbb{I}[f(x) \geq \tau] \mid f \in \mathcal{H}, \tau \in \mathbb{R} \right\}$$

$$= \sup_{h \in \mathcal{H}'} \left| \mathbb{P}_{\tilde{p}}[h(x) = 1] - \mathbb{P}_{\tilde{p}_n}[h(x) = 1] \right|$$

$$= \sup_{h \in \mathcal{H}'} \left| \sigma(h) - \sigma_n(h) \right| \leq \mathcal{O}\left( \sqrt{\frac{vc(\mathcal{H}') + \log(1/\delta)}{n}} \right)$$

$$\underbrace{\phantom{\sup_{h \in \mathcal{H}'} \left| \sigma(h) - \sigma_n(h) \right|}}_{\text{VC-dim notation}}$$

w.p $1-\delta$.

$$\widetilde{TV}(\tilde{p}, \tilde{p}_n) \doteq \sqrt{\frac{d}{n}}, \text{ where } d = vc(\mathcal{H}').$$
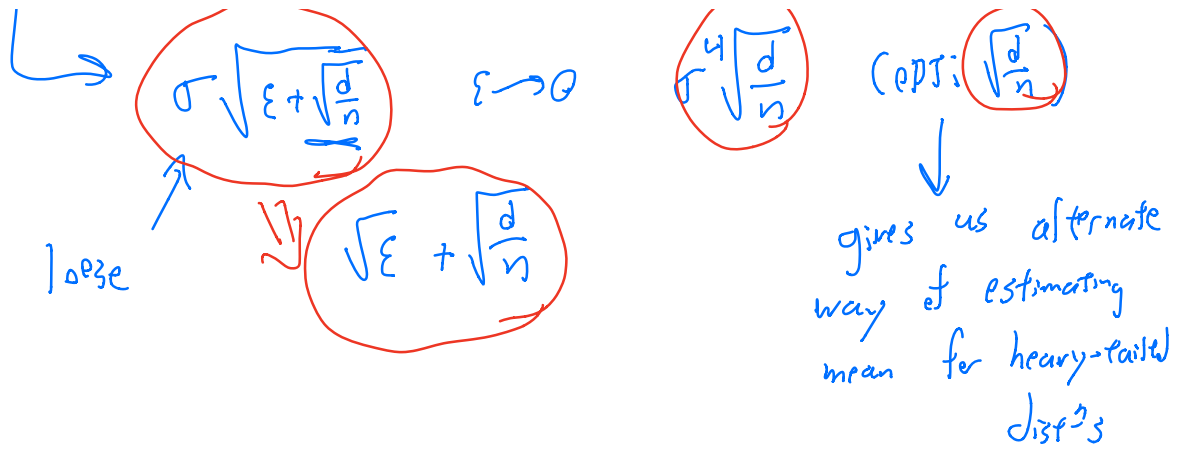
## Interpretation,

Resilience bound:

Bounded covariance: error $\mathcal{O}(\sigma\sqrt{\varepsilon})$

Subgaussian: error $\mathcal{O}\left( \sigma \varepsilon \sqrt{\log(1/\varepsilon)} \right)$.

$\varepsilon \mapsto \tilde{\varepsilon}$

$$\sigma \sqrt{\varepsilon + \sqrt{\frac{d}{n}}} \qquad \varepsilon \longrightarrow 0 \qquad \sqrt[4]{\frac{d}{n}} \qquad (\text{OPT}: \sqrt{\frac{d}{n}})$$

loose

$$\sqrt{\varepsilon} + \sqrt{\frac{d}{n}}$$

gives us alternate way of estimating mean for heavy-tailed dist's

Connections b/t $\widetilde{TV}$ and Tukey median