

Lec 7: Minimax lower bounds: Le Cam, Assouad, and Fano

Yanjun Han



Recall: P_θ : statistical model with unknown parameter $\theta \in \Theta$

$X \sim P_\theta$: observation

$\hat{\theta} = \hat{\theta}(x)$: estimator

$L(\theta, a) \geq 0$: loss function

Target of upcoming lectures: characterize the minimax risk

$$r^* = \inf_{\pi} \sup_{\theta \in \Theta} \mathbb{E}_\pi [L(\theta, \hat{\theta}(x))]$$

Upper bound: construct an estimator $\hat{\theta}$ and analyze its worst-case risk

Lower bound (*): use information-theoretic arguments to show a fundamental limit for any estimator

Last lecture (LAM): asymptotic, exact constant

This lecture: non-asymptotic, focus on the optimal rate

High-level idea: $r^* \geq r_\pi^*$ for any prior π

- constructing the least favorable prior and analyzing the Bayes risk can both be hard
- try simpler priors:
 - 1) binary hypothesis testing: $\pi = \text{Unif}(\{\theta_0, \theta_1\})$ (Le Cam's two point method)
 - 2) testing multiple hypotheses: $\pi = \text{Unif}(\{\theta_0, \dots, \theta_m\})$ (Assouad, Fano)

Le Cam's two point method

Thm. Suppose $\theta_0, \theta_1 \in \Theta$ satisfy the separation condition:

$$\inf_a (L(\theta_0, a) + L(\theta_1, a)) \geq \Delta.$$

Then

$$r^* \geq \inf_{\pi} \frac{1}{2} (\mathbb{E}_{\theta_0} L(\theta_0, \hat{\theta}) + \mathbb{E}_{\theta_1} L(\theta_1, \hat{\theta})) \geq \frac{\Delta}{2} (1 - \text{TV}(P_{\theta_0}, P_{\theta_1})).$$

General paradigm of applying two point method: find two points $\theta_0, \theta_1 \in \Theta$ satisfying

- separation condition: $\inf_{\alpha} (L(\theta_0, \alpha) + L(\theta_1, \alpha)) \geq \Delta$

(no single action performs uniformly well under θ_0, θ_1)

- indistinguishability condition: $TV(P_{\theta_0}, P_{\theta_1}) \leq 1 - \Delta L(1)$

(no test can reliably distinguish between P_{θ_0} and P_{θ_1})

By the joint range of TV vs. H^2, KL, χ^2 , the second condition is implied by:

$$\textcircled{1} \quad H^2(P_{\theta_0}, P_{\theta_1}) \leq 2 - \Delta L(1)$$

$$\textcircled{2} \quad D_{KL}(P_{\theta_0} \parallel P_{\theta_1}) = O(1) \quad \text{or} \quad D_{KL}(P_{\theta_1} \parallel P_{\theta_0}) = O(1)$$

$$\textcircled{3} \quad \chi^2(P_{\theta_0} \parallel P_{\theta_1}) = O(1) \quad \text{or} \quad \chi^2(P_{\theta_1} \parallel P_{\theta_0}) = O(1)$$

$$\textcircled{4} \quad I(\theta; X) \leq \log 2 - \Delta L(1), \text{ with } \theta \sim \text{Unif}(\{\theta_0, \theta_1\}) \quad (\text{exercise!})$$

Pf. For any estimator $\hat{\theta} = \hat{\theta}(X)$,

$$\begin{aligned} E_{\theta_0} L(\theta_0, \hat{\theta}) + E_{\theta_1} L(\theta_1, \hat{\theta}) &= \int L(\theta_0, \hat{\theta}(x)) P_{\theta_0}(x) dx + \int L(\theta_1, \hat{\theta}(x)) P_{\theta_1}(x) dx \\ &\geq \Delta \int \min\{P_{\theta_0}(x), P_{\theta_1}(x)\} dx \\ &= \Delta \int \frac{1}{2}(P_{\theta_0}(x) + P_{\theta_1}(x) - |P_{\theta_0}(x) - P_{\theta_1}(x)|) dx \\ &= \Delta (1 - TV(P_{\theta_0}, P_{\theta_1})) \end{aligned}$$

◻

Despite the simplicity of the two-point method, it has numerous applications.

Example 1.1 (normal mean estimation) $X \sim N(\theta, \sigma^2)$, unknown $\theta \in \mathbb{R}$, known σ^2 .

Target: $r^* = \inf_{\hat{\theta}} \sup_{\theta} E_{\theta}[(\hat{\theta} - \theta)^2]$.

Clearly, by choosing $\hat{\theta}(X) = X$, we achieve $r^* \leq \sigma^2$.

To lower bound r^* , apply two point method with $\theta_0 = 0, \theta_1 = \delta$.

- separation condition: $\Delta = \frac{\delta^2}{2}$

- indistinguishability condition: $TV(N(0, \sigma^2), N(\delta, \sigma^2)) = 2(1 - \Phi(\frac{|\delta|}{2\sigma}))$

Therefore, $r^* \geq \sup_{\delta \in \mathbb{R}} \frac{\delta^2}{2}(1 - \Phi(\frac{|\delta|}{2\sigma})) \approx 0.332 \sigma^2$.

(Compare with $r^* = \sigma^2$ by Anderson's lemma in Lec 6)

Example 1.2 (Binomial model) Let $X \sim B(n, \theta)$ with unknown $\theta \in [0, 1]$.

$$\text{Target: } r^* = \inf_{\hat{\theta}} \sup_{\theta} \mathbb{E}_{\theta}[(\hat{\theta} - \theta)^2]$$

Upper bound: choose $\hat{\theta}(X) = \frac{X}{n}$, then

$$\mathbb{E}_{\theta}\left[\left(\frac{X}{n} - \theta\right)^2\right] = \frac{\theta(1-\theta)}{n} \leq \frac{1}{4n} = O\left(\frac{1}{n}\right).$$

Lower bound: apply two point method to $\theta_0 = \frac{1}{2}$, $\theta_1 = \frac{1}{2} + \frac{1}{2\sqrt{n}}$

- separation: $\Delta = \frac{1}{2}\left(\frac{1}{2\sqrt{n}}\right)^2 = \Omega\left(\frac{1}{n}\right)$

- distinguishability: $D_{KL}(B(n, \theta_1) \parallel B(n, \theta_0)) = n \cdot D_{KL}(\text{Bern}(\theta_1) \parallel \text{Bern}(\theta_0))$

$$= \frac{n}{2} \left(\left(1 + \frac{1}{\sqrt{n}}\right) \log \left(1 + \frac{1}{\sqrt{n}}\right) + \left(1 - \frac{1}{\sqrt{n}}\right) \log \left(1 - \frac{1}{\sqrt{n}}\right) \right)$$

$$= \frac{n}{2} \cdot O\left(\left(\frac{1}{\sqrt{n}}\right)^2\right) = O(1).$$

This implies $r^* = \Omega\left(\frac{1}{n}\right)$ as well.

Example 1.3 (functional estimation) $X = (X_1, \dots, X_n)$ are i.i.d. draws from an unknown

pmf $P = (p_1, \dots, p_k)$. Loss: $L(P, \hat{P}) = |\hat{P} - H(P)|$ (entropy estimation)

Result (Jiao et al.'15; Wu and Yang '16):

$$r^* = \inf_{\hat{P}} \sup_P \mathbb{E}_P |\hat{P} - H(P)| \asymp \frac{k}{n \log n} + \frac{\log k}{\sqrt{n}} \quad \text{if } k \lesssim n \log n$$

Pf of the simpler $\Omega\left(\frac{\log k}{\sqrt{n}}\right)$ lower bound: since $D_{KL}(P_{X|P_0} \parallel P_{X|P_1}) = n D_{KL}(P_0 \parallel P_1)$

we solve the optimization program motivated by the two point method:

$$\max_{P_0} |H(P_0) - H(P_1)|$$

$$\text{s.t. } D_{KL}(P_0 \parallel P_1) \leq \frac{c}{n}$$

$$\text{Try } P_0 = \left(\frac{1}{2}, \frac{1}{2(k-1)}, \dots, \frac{1}{2(k-1)}\right), \quad P_1 = \left(\frac{1-\varepsilon}{2}, \frac{1+\varepsilon}{2(k-1)}, \dots, \frac{1+\varepsilon}{2(k-1)}\right)$$

$$\text{Then } D_{KL}(P_0 \parallel P_1) = \frac{1}{2} \log \frac{1}{1-\varepsilon^2} = O\left(\frac{1}{n}\right) \text{ if } \varepsilon = O\left(\frac{1}{\sqrt{n}}\right),$$

$$|H(P_0) - H(P_1)| = \left| \frac{1}{2} \log 2 + \frac{1}{2} \log 2(k-1) - \frac{1-\varepsilon}{2} \log \frac{2}{1-\varepsilon} - \frac{1+\varepsilon}{2} \log \frac{2(k-1)}{1+\varepsilon} \right| \\ \asymp \varepsilon \log k.$$

Therefore, choosing $\varepsilon \asymp \frac{1}{\sqrt{n}}$ proves that $r^* = \Omega\left(\frac{\log k}{\sqrt{n}}\right)$. □

(The other lower bound $\Omega\left(\frac{k}{n \log n}\right)$ is proven using a more involved two point method, which will be the topic of Lec 8)

Example 1.4 (two-armed bandit) $\theta = (\mu_1, \mu_2) \in [0, 1]^2$

Observation: for $t \in [T]$, learner pulls an arm $\pi_t \in \{1, 2\}$ based on (π^{t-1}, r^{t-1}) .
and observes reward $r_t \sim N(\mu_{\pi_t}, 1)$.

Learner aims to minimize the regret $R_T(\pi) = T \max\{\mu_1, \mu_2\} - \sum_{t=1}^T \mu_{\pi_t}$.

We will show that $r^* = \inf_{\pi} \sup_{\substack{\mu_1, \mu_2 \\ |\mu_1 - \mu_2| > \Delta}} \mathbb{E}_{\mu_1, \mu_2} [R_T(\pi)] = \Omega\left(\frac{1 \vee \log(T\Delta^2)}{\Delta} \wedge T\Delta\right)$

(In particular, choosing $\Delta \asymp \frac{1}{\sqrt{T}}$ gives the usual lower bound $\Omega(\sqrt{T})$ for two-armed bandit)

Pf. First, by the chain rule of KL, we have (exercise)

$$\begin{aligned} D_{KL}(P_{\mu_1, \mu_2} \| P_{\mu'_1, \mu'_2}) &= \sum_{t=1}^T \mathbb{E}_{P_{\mu_1, \mu_2}} \left[\frac{(\mu_1 - \mu'_1)^2}{2} \mathbb{1}(\pi_t=1) + \frac{(\mu_2 - \mu'_2)^2}{2} \mathbb{1}(\pi_t=2) \right] \\ &= \frac{(\mu_1 - \mu'_1)^2}{2} \mathbb{E}_{P_{\mu_1, \mu_2}} [T_1] + \frac{(\mu_2 - \mu'_2)^2}{2} \mathbb{E}_{P_{\mu_1, \mu_2}} [T_2], \end{aligned}$$

where $T_i = \sum_{t=1}^T \mathbb{1}(\pi_t=i)$ is the total number of times we pull arm $i \in \{1, 2\}$.

Motivated by this, choose two points $(\mu_1, \mu_2) = (\Delta, 0)$, $(\mu'_1, \mu'_2) = (\Delta, 2\Delta)$.

The separation parameter is $T\Delta$ (why?). and

$$D_{KL}(P_{\mu_1, \mu_2} \| P_{\mu'_1, \mu'_2}) = 2\Delta^2 \mathbb{E}_i[T_i] \quad (\mathbb{E}_i := \mathbb{E}_{P_{\mu_1, \mu_2}})$$

Two point method then gives $r^* = \Omega(T\Delta \cdot \exp(-2\Delta^2 \mathbb{E}_i[T_i]))$

$$(1 - TV(P, Q)) \geq \frac{1}{2} \exp(-D_{KL}(P \| Q))$$

Note that $\mathbb{E}_i[T_i]$ depends on the policy π , and this lower bound is useful only if $\mathbb{E}_i[T_i]$ is small. To address it, note that we have a different lower bound:

$$r^* \geq \mathbb{E}_i[R_T(\pi)] = \Delta \cdot \mathbb{E}_i[T_2].$$

thus

$$\begin{aligned} r^* &= \Omega(\max\{\Delta \cdot \mathbb{E}_i[T_2], T\Delta \cdot \exp(-2\Delta^2 \mathbb{E}_i[T_2])\}) \\ &= \Omega(\min_{T_0 \in [0, T]} \max\{\Delta T_0, T\Delta \cdot \exp(-2\Delta^2 T_0)\}) \\ &= \Omega\left(\frac{1 \vee \log(T\Delta^2)}{\Delta} \wedge T\Delta\right) \end{aligned}$$

□

Example 1.5 (multi-armed bandit) Same observation model, but now with K arms:
 $\theta = (\mu_1, \dots, \mu_K) \in [0, 1]^K$, $R_T(\pi) = T \max_{i \in [K]} \mu_i - \sum_{t=1}^T \mu_{\pi_t}$.

We will show that $r^* = \inf_{\pi} \sup_{\theta} \mathbb{E}_{\theta} [R_T(\pi)] = \Omega(\sqrt{KT})$.

(Interestingly, two points suffice for this example!)

Pf. Choose $\theta_1 = (\Delta, 0, 0, \dots, 0)$

$$\theta_{2,i} = (\Delta, 0, \dots, 0, 2\Delta, 0, \dots, 0), \quad i = 2, \dots, K.$$

For each pair $(\theta_1, \theta_{2,i})$, the separation parameter is always $T\Delta$.

In addition, for any policy π , $D_{KL}(P_{\theta_1} \| P_{\theta_{2,i}}) = 2\Delta^2 \mathbb{E}_i[T_i]$.

$$(\text{recall } T_i = \sum_{t=1}^T \mathbf{1}(\pi_t = i))$$

Key observation: since $\sum_{i=2}^K \mathbb{E}_i[T_i] \leq T$, there must $\exists i_0$ s.t. $\mathbb{E}_{i_0}[T_{i_0}] \leq \frac{T}{K-1}$.

Now applying two point argument to $(\theta_1, \theta_{2,i_0})$ and $\Delta \asymp \sqrt{\frac{K}{T}}$ gives

$$D_{KL}(P_{\theta_1} \| P_{\theta_{2,i_0}}) = O(1), \text{ and therefore}$$

$$r^* = \Omega(T\Delta) = \Omega(\sqrt{KT}).$$

④

Testing multiple hypotheses

Why two points may fail? Look at normal mean estimation in high dimensions:

$$X \sim N(\theta, \sigma^2 I_n), \text{ with loss } L(\theta, \hat{\theta}) = \|\hat{\theta} - \theta\|_2^2.$$

Two point method gives at best

$$r^* \geq \sup_{\theta_0, \theta_1} \frac{\|\theta_0 - \theta_1\|_2^2}{2} \left(1 - \Phi\left(\frac{\|\theta_0 - \theta_1\|_2}{2\sigma}\right)\right) \asymp \sigma^2.$$

This doesn't capture the dependence on the dimension n ! (Recall that $r^* = n\sigma^2$ by Anderson's lemma)

At a high level, this is because testing between two hypotheses does NOT capture the true difficulty of a high-dimensional problem!

Challenge in high dimensions:

- separation condition: there might be different separation conditions
- indistinguishability condition: instead of the binary testing case where $1 - TV(P, Q)$ is a tight measure of the test measure, this tight measure in the multiple hypotheses case is often not tractable (and needs further lower bound)

Pairwise separation: Fano's inequality

Thm. Let $\theta_1, \dots, \theta_m \in \Theta$ satisfy

$$\min_{i \neq j} \inf_a (L(\theta_i, a) + L(\theta_j, a)) \geq \Delta.$$

Then for $\pi = \text{Unif}(\{\theta_1, \dots, \theta_m\})$,

$$r_\pi^* \geq \frac{\Delta}{2} \left(1 - \frac{I(\theta; X) + \log 2}{\log m} \right), \quad \theta \sim \pi, X | \theta \sim P_\theta.$$

Before we prove it, we establish a useful "golden formula" for mutual info:

$$\text{Lemma. } I(X; Y) = \min_{Q_Y} D_{KL}(P_{XY} \| P_X Q_Y) = \min_{Q_Y} \mathbb{E}_{P_X} [D_{KL}(P_{Y|X} \| Q_Y)]$$

Pf. Simply note that $I(X; Y) = D_{KL}(P_{XY} \| P_X Q_Y) - D_{KL}(P_Y \| Q_Y), \forall Q_Y$. \blacksquare

Pf of Fano. $P_{\theta X}$

$$(P_{\theta X}) \xrightarrow{(P_{\theta} P_X)} \xrightarrow{(P_{\theta} P_X)} \text{Bern}(P(L(\theta, \hat{\theta}) < \frac{\Delta}{2})) \xrightarrow{\text{Bern}(q)}$$

Note that $q \leq \frac{1}{m}$ by separation condition. Then DPI gives

$$D_{KL}(\text{Bern}(P(L(\theta, \hat{\theta}) < \frac{\Delta}{2})) \| \text{Bern}(q)) \leq I(\theta; X)$$

$$\Rightarrow P(L(\theta, \hat{\theta}) \geq \frac{\Delta}{2}) \geq 1 - \frac{I(\theta; X) + \log 2}{\log m}. \text{ The rest follows from Markov's ineq. } \blacksquare$$

In fact, the previous argument establishes a general result:

Thm (generalized Fano). For any prior π and $\Delta > 0$,

$$r_{\pi}^* \geq \Delta \left(1 - \frac{I(\theta; X) + \log 2}{\log \left(\frac{1}{P_{\Delta}} \right)} \right), \quad \theta \sim \pi, \quad X|\theta \sim P_{\theta}$$

where $P_{\Delta} = \sup_a \pi(L(\theta, a) < \Delta)$ is the "small-ball probability".

The classical Fano's inequality is a special case with $\pi \sim \text{Unif}(\{\theta_1, \dots, \theta_m\})$ and a pairwise separation condition.

Additive separation: Assouad's lemma

Thm. For a hypercube parametrization $u \in \{\pm 1\}^d$, associate $\theta_u \in \mathbb{R}$.

Suppose $\inf_a (L(\theta_u, a) + L(\theta_{u'}, a)) \geq \Delta - \sum_{i=1}^d 1(u_i \neq u'_i)$

then for the prior $\pi = \text{Unif}(\{\theta_u\}_{u \in \{\pm 1\}^d})$,

$$r_{\pi}^* \geq \frac{\Delta}{4} \sum_{i=1}^d (1 - \text{TV}(P_{i,+}, P_{i,-})),$$

where $P_{i,+} = \frac{1}{2^{d-1}} \sum_{u: u_i=1} P_{\theta_u}, \quad P_{i,-} = \frac{1}{2^{d-1}} \sum_{u: u_i=-1} P_{\theta_u}.$

The following corollaries are often used:

Corollary 1: $r_{\pi}^* \geq \frac{\Delta}{4} (1 - \max_{u, u' \text{ neighbors}} \text{TV}(P_{\theta_u}, P_{\theta_{u'}})).$ (the classical Assouad)

Corollary 2: $r_{\pi}^* \geq \frac{\Delta}{4} (1 - \mathbb{E}_{u \sim \text{Unif}(\{\pm 1\}^d)} \mathbb{E}_{i \sim \text{Unif}([d])} \text{TV}(P_{\theta_u}, P_{\theta_{u \oplus i}}))$
 $u \text{ with } i\text{-th bit flipped}$

Pf of Assouad. For any estimator $\hat{\theta}$, construct $\hat{u} = (\hat{u}_1, \dots, \hat{u}_d) \in \{\pm 1\}^d$ as follows:

$$\hat{u} = \arg_{\hat{u}} \min L(\theta_u, \hat{\theta}).$$

Then $\forall u \in \{\pm 1\}^d$,

$$\begin{aligned} L(\theta_u, \hat{\theta}) &\geq \frac{L(\theta_u, \hat{\theta}) + L(\theta_{\hat{u}}, \hat{\theta})}{2} \geq \frac{\Delta}{2} \sum_{i=1}^d I(u_i \neq \hat{u}_i). \\ &\Rightarrow \frac{1}{2^d} \sum_u \mathbb{E}_{\theta_u} [L(\theta_u, \hat{\theta})] \geq \frac{1}{2^d} \sum_u \frac{\Delta}{2} \sum_{i=1}^d P_{\theta_u}(u_i \neq \hat{u}_i) \\ &= \frac{\Delta}{4} \sum_{i=1}^d (P_{i,+}(\hat{u}_i \neq 1) + P_{i,-}(\hat{u}_i \neq -1)) \\ &\geq \frac{\Delta}{4} \sum_{i=1}^d (1 - TV(P_{i,+}, P_{i,-})). \quad \blacksquare \end{aligned}$$

(Exercise: show that $\frac{1}{d} \sum_{i=1}^d TV(P_{i,+}, P_{i,-}) = 1 - \Omega(1) \iff \frac{1}{d} \sum_{i=1}^d I(U_i; X) = \log 2 - \Omega(1)$
for $U \sim \text{Unif}(\{\pm 1\}^d)$.

In addition, since $I(U; X) \geq \sum_{i=1}^d I(U_i; X)$, under this hypercube condition, Assouad is no weaker than Fano.)

Fano and Assouad are good at proving lower bounds for high-dimensional targets.

Example 2.1 (normal mean model) $X \sim N(\theta, \sigma^2 I_n)$, $L(\theta, \hat{\theta}) = \|\theta - \hat{\theta}\|_2^2$.

We will show that $r^* = \Omega(n\sigma^2)$.

Pf 1 (Fano). Construct a subset \mathbb{B}_0 of $\{\pm \delta\}^n$ (with δ TBD) s.t.

- $m := |\mathbb{B}_0|$ is large enough;
- $\min_{\theta \neq \theta' \in \mathbb{B}_0} \|\theta - \theta'\|_2^2 \geq \frac{\delta^2 n}{5}$.

By Gilbert-Varshamov bound below, we can make $m = e^{\Omega(n)}$. Then $\Delta = \frac{\delta^2 n}{10}$, and

$$I(\theta; X) \leq \max_{\substack{\theta \in \mathbb{B}_0}} D_{KL}(N(\theta, \sigma^2 I_n) \| N(0, \sigma^2 I_n)) = \frac{n\delta^2}{2\sigma^2}.$$

$$\text{Fano} \Rightarrow r^* = \Omega\left(\delta^2 n \left(1 - \frac{n\delta^2 \log 2}{\Omega(n)}\right)\right) \xrightarrow{\delta \asymp \sigma} r^* = \Omega(n\sigma^2).$$

Lemma (Gilbert-Varshamov) $\exists A \subseteq \{\pm 1\}^n$ s.t. $\min_{\substack{u, u' \in A \\ u \neq u'}} \sum_{i=1}^n \mathbb{1}(u_i \neq u'_i) \geq d$, and

$$m \geq \frac{2^n}{\sum_{j=0}^{d-1} \binom{n}{j}} = 2^{n(1 - h_2(\frac{1}{2}) + o(n))} \quad (h_2(x) = x \log_2 \frac{1}{x} + (1-x) \log_2 \frac{1}{1-x})$$

Pf. Volume argument: $\forall u \in \{\pm 1\}^n$, $|\{u' \in \{\pm 1\}^n : \sum_{i=1}^n \mathbb{1}(u_i \neq u'_i) \leq d-1\}| = \sum_{j=0}^{d-1} \binom{n}{j}$.
 So if $m < \frac{2^n}{\sum_{j=0}^{d-1} \binom{n}{j}}$, there must exist some $u \in \{\pm 1\}^n$ with distance $\geq d$ to all existing points. \square

Pf 2 (Generalized Fano) Let $\theta \sim \text{Unif}(\{\pm \delta\}^n)$, then $I(\theta; X) \leq \frac{n\delta^2}{2\sigma^2}$
 Pick $\Delta = \frac{n\delta^2}{12}$, then the small-ball probability is

$$p_0 = \sup_a \pi(\|\theta - a\|_2^2 < \Delta) \leq \sup_{\substack{\theta \in \{\pm \delta\}^n \\ \hat{\theta} = \arg \min \|\theta - \hat{\theta}\|_2^2}} \pi(\|\theta - \hat{\theta}\|_2^2 < 4\Delta) = \frac{1}{2^n} \sum_{j=0}^{\lceil n/3 \rceil - 1} \binom{n}{j} = 2^{-\Omega(n)}$$

↑ Stirling

Therefore, we again get $r^* = \Omega(n\delta^2(1 - \frac{n\delta^2 + 4\Delta}{\Omega(n)})) \stackrel{\delta \approx \sigma}{=} \Omega(n\sigma^2)$.

Pf 3 (Assouad). For $\delta > 0$, let $\theta_u = \delta u$, $u \in \{\pm 1\}^n$. Then $\Delta = 2\delta^2$, and

$$\max_{u, u' \text{ neighbors}} \text{TV}(N(\theta_u, \sigma^2 I_n), N(\theta_{u'}, \sigma^2 I_n)) = 2(1 - \Phi(\frac{\delta}{\sigma})).$$

Choosing $\delta = \sigma$, Assouad's lemma gives $r^* = \Omega(n\Delta) = \Omega(n\sigma^2)$.

Example 2.2 (learning theory). $(X_1, Y_1), \dots, (X_n, Y_n) \sim P_{XY}$ (unknown) with $Y_i \in \{0, 1\}$.

Let F be a given class of functions $X \rightarrow \{0, 1\}$ with VC dimension d .

Excess risk for a trained classifier $\hat{f}: X \rightarrow \{0, 1\}$ based on $(X_1, Y_1), \dots, (X_n, Y_n)$:

$$ER(\hat{f}) := ERL(\hat{f}; P_{XY}) := P_{XY}[Y \neq \hat{f}(X)] - \min_{f \in F} P_{XY}[Y \neq f(X)]$$

We will show that for $n \geq d$,

$$\inf_{\hat{f}} \sup_{P_{XY}} E_{P_{XY}} [ER(\hat{f})] = \Omega(\sqrt{\frac{d}{n}}) \quad (\text{agnostic setting})$$

$$\inf_{\hat{f}} \sup_{\substack{P_{XY}: \exists f \in F \\ Y=f(X) \text{ P}_{XY}-\text{a.s.}}} E_{P_{XY}} [ER(\hat{f})] = \Omega(\frac{d}{n}) \quad (\text{realizable setting})$$

Recall: $\text{VC-dim}(\mathcal{F}) = d \Rightarrow \exists x_1, \dots, x_d \in X, \forall u \in \{\pm 1\}^d, \exists f_u \in \mathcal{F} \text{ s.t. } f_u(x_i) = u_i \forall i \in [d]$.

Pf (agnostic case). Use $x_1, \dots, x_d \in X$ and $\{f_u\}_{u \in \{\pm 1\}^d} \subseteq \mathcal{F}$ in the above definition.

Under $u \in \{\pm 1\}^d$, construct $P_u := P_{XY, u}$ as:

$$X \sim \text{Unif}(\{x_1, \dots, x_d\}), \quad Y | X=x_i = \begin{cases} u_i & \text{w.p. } \frac{1}{2} + \delta \\ -u_i & \text{w.p. } \frac{1}{2} - \delta \end{cases}$$

- Separation condition: $\min_{f \in \mathcal{F}} P_u(f(X) \neq Y) = \frac{1}{2} - \delta, \forall u$.

$$\forall f_0: \text{ER}(f_0, P_u) + \text{ER}(f_0, P_w)$$

$$= P_u(Y \neq f_0(X)) + P_w(Y \neq f_0(X)) - 2\left(\frac{1}{2} - \delta\right)$$

$$\geq \sum_{i=1}^d \frac{1}{d} (1(u_i \neq u'_i) \cdot 1 + 1(u_i = u'_i) \cdot (1-2\delta)) - (1-2\delta)$$

$$= \frac{2\delta}{d} \sum_{i=1}^d 1(u_i \neq u'_i) \Rightarrow \Delta = \frac{2\delta}{d}.$$

- Indistinguishability condition: for neighboring u, u' :

$$\begin{aligned} D_{KL}(P_u^{\otimes n} \| P_{u'}^{\otimes n}) &= n D_{KL}(P_u \| P_{u'}) = n \cdot \frac{1}{d} D_{KL}(\text{Bern}(\frac{1}{2} + \delta) \| \text{Bern}(\frac{1}{2} - \delta)) \\ &= O\left(\frac{n\delta^2}{d}\right) \text{ for } \delta = \frac{1}{2} - \Omega(1). \end{aligned}$$

Finally, choosing $\delta \approx \sqrt{\frac{d}{n}}$. Assouad's lemma gives $r^* = \Omega(d\Delta) = \Omega(\sqrt{\frac{d}{n}})$.

(Realizable case) For $u \in \{\pm 1\}^d$, now define $P_u := P_{XY, u}$ as:

$$X = \begin{cases} x_i & \text{w.p. } 1 - \frac{d-1}{n} \\ x_i & \text{w.p. } \frac{1}{n} \quad \forall 2 \leq i \leq d \end{cases}, \quad Y | X=x_i = u_i \text{ w.p. 1.}$$

Clearly $\min_{f \in \mathcal{F}} P_u(f(X) \neq Y) = 0, \forall u$.

- Separation condition: $\Delta = \frac{1}{n}$ by a similar analysis

- Indistinguishability condition: for $u' = u^{\otimes i}$,

$$TV(P_u^{\otimes n}, P_{u'}^{\otimes n}) \leq P(x_i \text{ appears in } X_1, \dots, X_n) = 1 - \left(1 - \frac{1}{n}\right)^n = 1 - \Omega(1).$$

Therefore. Assouad's lemma gives $r^* = \Omega((d-1)\Delta) = \Omega(\frac{d}{n})$.

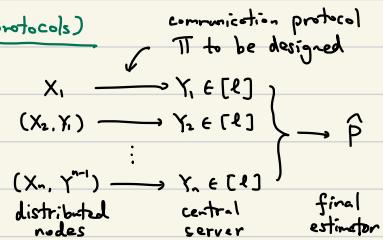
(See HW for a further generalization.)

Assouad's lemma is also surprisingly flexible in sequential settings.

Example 2.3 (distribution estimation under sequential communication protocols)

Let $P = (P_1, \dots, P_k)$ be an unknown pmf.

The observations are $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} P$, but \hat{P} must be formed via the distributed diagram on the right, together with communication constraints $\ell \leq k$.



We will show that, $r^* = \inf_{(\beta, \Pi)} \sup_P \mathbb{E}_P [\text{TV}(P, \beta)] = \Omega\left(\frac{k}{\sqrt{n\ell}}\right)$ if $k \leq \ell$
 $n \geq \frac{k^2}{\ell}$.

Pf. WLOG assume k is even. For $u \in \{\pm 1\}^{k/2}$, construct

$$P_u = \left(\frac{1+\delta u_1}{k}, \frac{1-\delta u_1}{k}, \dots, \frac{1+\delta u_{k/2}}{k}, \frac{1-\delta u_{k/2}}{k} \right), \quad \delta \in (0, \frac{1}{2}) \text{ TBD.}$$

Easy to check that $\Delta = \Omega\left(\frac{\delta}{k}\right)$. Next we use Corollary 2 of Assouad's lemma and try to upper bound

$$\begin{aligned} \mathbb{E}_u \mathbb{E}_i [\text{TV}(P_{Y^i|u}, P_{Y^i|u^\perp})] &\leq \sqrt{\mathbb{E}_u \mathbb{E}_i [\text{TV}(P_{Y^i|u}, P_{Y^i|u^\perp})^2]} \quad (\text{Jensen}) \\ &\leq \sqrt{\mathbb{E}_u \mathbb{E}_i [H^2(P_{Y^i|u}, P_{Y^i|u^\perp})]} \quad (\text{TV} \leq H) \\ &\leq C \sqrt{\sum_{t=1}^{\infty} \mathbb{E}_u \mathbb{E}_i [H^2(P_{Y_t|Y^{t-1}, u}, P_{Y_t|Y^{t-1}, u^\perp})]} \quad (\text{Jayaram's subadditivity of } H^2 \text{ in Lec 3}) \end{aligned}$$

Note that $P_{Y_t|Y^{t-1}, u} = \sum_{x \in [k]} P_{Y_t=Y^t, X_t=x} \cdot P_u(x) \geq \frac{1}{2k} \sum_{x \in [k]} P_{Y_t=Y^t, X_t=x}$,

so

$$\begin{aligned} x^2(P_{Y_t|Y^{t-1}, u^\perp} || P_{Y_t|Y^{t-1}, u}) &\leq \sum_{y \in [k]} \frac{(P_{Y_t=Y^t, u^\perp} - P_{Y_t=Y^t, u})^2}{\frac{1}{2k} \sum_{x \in [k]} P_{Y_t=y | Y^{t-1}, X_t=x}} \\ &= 2k \left(\frac{2\delta}{k}\right)^2 \sum_{y \in [k]} \frac{(P_{Y_t=Y^t, X_t=2i-1} - P_{Y_t=Y^t, X_t=2i})^2}{\sum_{x \in [k]} P_{Y_t=y | Y^{t-1}, X_t=x}} \\ &\leq \frac{8\delta^2}{k} \sum_{y \in [k]} \frac{P_{Y_t=Y^t, X_t=2i-1} + P_{Y_t=Y^t, X_t=2i}}{\sum_{x \in [k]} P_{Y_t=y | Y^{t-1}, X_t=x}}. \end{aligned}$$

Therefore,

$$\mathbb{E}_i[H^2(P_{T+1|Y^{t+1}, u}, P_{T+1|Y^{t+1}, u^*})] \leq \mathbb{E}_i[X^2(P_{T+1|Y^{t+1}, u^*} \| P_{T+1|Y^{t+1}, u})] \\ = O\left(\frac{\delta^2 \ell}{k^2}\right).$$

Putting everything together, Assouad's lemma yields

$$r^* = \Omega\left(\delta(1 - O\left(\sqrt{\frac{n\delta^2 \ell}{k^2}}\right))\right) \stackrel{\delta \approx \frac{k}{\sqrt{n\ell}}, \delta < \frac{1}{2}}{=} \Omega\left(\frac{k}{\sqrt{n\ell}}\right).$$

(4)

Special topic : interactive Le Cam

Model for interactive decision making :

- unknown true model M^* in a given model class \mathcal{M} (e.g. reward distribution of all arms)
- at each round $t = 1, \dots, T$:
 - ① learner chooses an action $a_t \in A$,
 - ② nature reveals reward $r_t \in [0, 1]$ and possible additional observation a_t , with $\mathbb{E}[r_t | a_t = a] = r^{M^*}(a)$, and $(r_t, a_t) \sim M^*(a_t)$
- learner aims to minimize the regret

$$R_T = \sum_{t=1}^T (r_*^{M^*} - r^{M^*}(a_t)), \quad \text{where } r_*^M = \max_{a \in A} r^M(a).$$

Example (multi-armed bandit)

- $A = [k]$
- $M^r = \{ \text{Bern}(\mu_i) \}_{i=1}^k, \quad \mathcal{M} = \{ M^{\mu} : \mu \in [0, 1]^k \}.$
- $M^r(a) = \text{Bern}(\mu_a), \quad r^{M^r}(a) = \mu_a.$

Question: What is a general two point lower bound for R_T ?

Idea: let $g^M(a) := r_*^M - r^M(a)$ denote the "gap" of $a \in A$ under M , then point lower bound suggests that

$$\inf_{\{a_t\}} \sup_{M^*} \mathbb{E}[R_T] \geq T \cdot \begin{cases} \sup_{M_0, M_1 \in \mathcal{M}} (\inf_{a \in A} g^{M_0}(a) + g^{M_1}(a)) \\ \text{s.t. } H^2(M_0, M_1) \leq \frac{c}{T} \text{ for small } c > 0. \end{cases}$$

- Challenges:
- ① The metric $\inf_a (g^{M_0}(a) + g^{M_1}(a))$ could be too pessimistic. For example, if a policy uses an action distribution p on a under M_0 , then $\mathbb{E}_p[g^{M_1}(a)]$ might be a better separation metric.
 - ② $H^2(M_0, M_1)$ is NOT well-defined, as the distributions of (r_t, a_t) depend on a_t . i.e. $(r_t, a_t) \sim M^*(a_t)$. So $H^2(M_0, M_1)$ should be replaced by $\mathbb{E}_p[H^2(M_0(a), M_1(a))]$ for some p , and take \inf_p somewhere.
 - ③ Where to take \inf_p (learner as the min player)?
 - $\sup_{M_0, M_1} \inf_p$: too small, by the same reason in ①
 - $\inf_p \sup_{M_0, M_1}$: too large, as the learner can adjust the action distribution in the sequential setting.

Defn. The constrained decision-to-estimation coefficient (DEC) is defined as

$$\text{dec}_{\varepsilon}(M) = \sup_{\bar{M}} \inf_{p \in \Delta(A)} \sup_{M \in M \cup \{\bar{M}\}} \left\{ \mathbb{E}_p[g^M(a)] : \mathbb{E}_p[H^2(M(a), \bar{M}(a))] \leq \varepsilon^2 \right\}.$$

- This is a sup-inf-sup structure: first choose a reference model \bar{M} , learner chooses an action distribution p based on \bar{M} , finally nature chooses an alternative model M
- The separation condition is w.r.t. \mathbb{E}_p
- The reference model \bar{M} doesn't need to belong to M

Example 3.1 (two-armed bandit). For two-armed bandit $(\text{Bern}(\mu_1), \text{Bern}(\mu_2))$ with $|\mu_1 - \mu_2| \geq \alpha$, choose $\bar{M} = (\text{Bern}(\frac{1}{2} + \delta), \text{Bern}(\frac{1}{2}))$ and $M \in \{(\text{Bern}(\frac{1}{2} + \delta), \text{Bern}(\frac{1}{2})), (\text{Bern}(\frac{1}{2} + \delta), \text{Bern}(\frac{1}{2} + \delta + \delta))\}$.

$$\text{dec}_{\varepsilon}(M) \geq \inf_{p \in [0,1]} \max \left\{ p_1 \delta, (1-p_1) \left(\frac{c\varepsilon}{\sqrt{p_1}} - \delta \right)_+ \right\} = \Omega(\delta \wedge \frac{\varepsilon^2}{\delta}).$$

Example 3.2 (multi-armed bandit) For $M = \left\{ \sum_{i=1}^K \text{Bern}(\mu_i) : \mu_1, \dots, \mu_K \in [0,1] \right\}$,

can choose $\bar{M} = \sum_{i=1}^K \text{Bern}(\frac{1}{2})$, pick $i_0 = \arg \min_i p_i$ and set $M(i_0) = \text{Bern}(\frac{1}{2} + \varepsilon \sqrt{K})$ to get $\text{dec}_{\varepsilon}(M) = \Omega(\varepsilon \sqrt{K})$ if $\varepsilon \sqrt{K} = O(1)$.

Thm (DFC lower bound). There exist absolute constants $c, C > 0$ s.t.

$$\inf_{\{a_t\}} \sup_{M^* \in M} \mathbb{E}_{M^*}[R_T] = \Omega\left(T \cdot (\text{dec}_\varepsilon(M) - C\varepsilon)_+\right), \text{ for } \varepsilon = \sqrt{\frac{c}{T}}.$$

Specializing to previous examples, this gives a lower bound $\Omega(\frac{1}{\Delta})$ for Example 3.1 when $\Delta \geq \frac{C}{T}$, and $\Omega(\sqrt{KT})$ for $T \geq K$.

Pf. (Simpler case: $\bar{M} \in M$, from (Foster, Golowich, Hsu '23)) Let $\Delta := \text{dec}_\varepsilon(M)$.

$$P_{\bar{M}} := \mathbb{E}_{\bar{M}}\left[\frac{1}{T} \sum_{t=1}^T p_t(\cdot | H^{t-1})\right]: \text{learner's average play under } \bar{M}$$

M : inner maximizer under $p = P_{\bar{M}}$

$$P_M := \mathbb{E}_M\left[\frac{1}{T} \sum_{t=1}^T p_t(\cdot | H^{t-1})\right]: \text{learner's average play under } M.$$

By definition,

$$\mathbb{E}_{a \sim P_{\bar{M}}} [g^M(a)] \geq \Delta \quad (1)$$

$$\mathbb{E}_{a \sim P_{\bar{M}}} [H^2(M(a), \bar{M}(a))] \leq \varepsilon^2 \quad (2)$$

By the sake of contradiction, we can assume that

$$\mathbb{E}_{a \sim P_M} [g^M(a)] \leq \frac{\Delta}{100} \quad (3)$$

$$\mathbb{E}_{a \sim P_{\bar{M}}} [g^{\bar{M}}(a)] \leq \frac{\Delta}{100} \quad (4)$$

We introduce some useful lemmas.

Lemma 1. For $c > 0$ small enough. $TV(P_M, P_{\bar{M}}) \leq 0.1$.

$$\begin{aligned} TV(P_M, P_{\bar{M}})^2 &\leq TV(P_{M,0T}^M, P_{M,0T}^{\bar{M}})^2 \quad (\text{DPI}) \\ &\leq H^2(P_{M,0T}^M, P_{M,0T}^{\bar{M}}) \\ &\leq C \sum_{t=1}^T \mathbb{E}_{a \sim P_{\bar{M}}} [H^2(P_{M,0t+1|H^{t-1}}^M, P_{M,0t+1|H^{t-1}}^{\bar{M}})] \quad (\text{subadditivity of } H^2) \\ &= C \sum_{t=1}^T \mathbb{E}_{a \sim P_{\bar{M}}} [H^2(M(a_t), \bar{M}(a_t))] \\ &= CT \cdot \mathbb{E}_{a \sim P_{\bar{M}}} [H^2(M(a), \bar{M}(a))] \\ &\leq 0.1 \quad (\text{by (2)}) \end{aligned} \quad (5)$$

Lemma 2. $\mathbb{E}_{\alpha \in \mathcal{A}} |r^M(\alpha) - r^{\bar{M}}(\alpha)| \leq \varepsilon$. (This step critically uses that the rewards are observed)

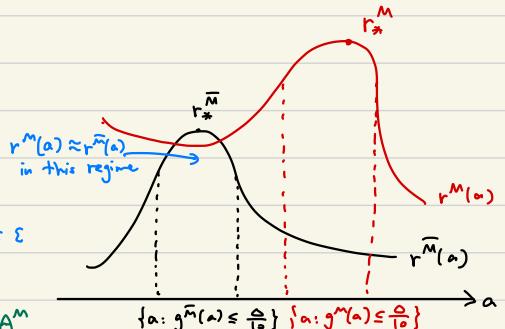
Pf. As $r_t \in [0, 1]$,

$$\text{LHS} \leq \mathbb{E}_{\alpha \in \mathcal{A}} [\text{TV}(M(\alpha), \bar{M}(\alpha))] \leq \mathbb{E}_{\alpha \in \mathcal{A}} [H(M(\alpha), \bar{M}(\alpha))] \leq \varepsilon. \quad \square$$

Next we present the proof.

1) By Lemma 2 and ①:

$$\begin{aligned} \Delta &\leq r_*^M - \mathbb{E}_{\alpha \in \mathcal{A}} [r^M(\alpha)] \\ &\leq r_*^M - \mathbb{E}_{\alpha \in \mathcal{A}} [r^{\bar{M}}(\alpha)] + \varepsilon \\ &= r_*^M - r_*^{\bar{M}} + \mathbb{E}_{\alpha \in \mathcal{A}} [g^{\bar{M}}(\alpha)] + \varepsilon \\ &\stackrel{④}{\leq} r_*^M - r_*^{\bar{M}} + \frac{\Delta}{100} + \varepsilon \\ \Rightarrow r_*^M - r_*^{\bar{M}} &\geq \frac{99\Delta}{100} - \varepsilon. \end{aligned}$$



2) By ③ and Markov: $P_M(\{a : g^M(a) \leq \frac{\Delta}{100}\}) \geq \frac{9}{10}$

3) By 2) and Lemma 1: $P_{\bar{M}}(A^M) \geq \frac{4}{5}$

4) By 1) and 3): $\mathbb{E}_{\alpha \in \mathcal{A}} [(r^M(\alpha) - r^{\bar{M}}(\alpha)) \mathbf{1}_{\alpha \in A^M}] \geq (r_*^M - \frac{\Delta}{100} - r_*^{\bar{M}}) \cdot P_{\bar{M}}(A^M) \geq (\frac{89\Delta}{100} - \varepsilon) \cdot \frac{4}{5}$

However, Lemma 2 states that LHS $\leq \varepsilon$. This is a contradiction when $\Delta > C\varepsilon$!

(General \bar{M} , from (Glasgow and Rakhlin '23))

For $\bar{M} \neq M$, ④ is no longer a result of small regret.

Solution: a stopping time argument. Let ALG be the original learner's algorithm, define

ALG' as follows: $\text{ALG}'_t = \text{ALG}_t$ as long as $\sum_{s=t}^T g^{\bar{M}}(a_s) < \frac{\Delta T}{2}$,
and ALG'_t always pulls $a^* = \arg \max_a r^{\bar{M}}(a)$ o.w.

Now redefine $P_{\bar{M}}, M, P_M$ using ALG' , then ①②④ + Lemma 1 & 2 still hold.

Let $\tau > 0$ be the stopping time of $\sum_{t=1}^{\tau} g^{\bar{M}}(a_t) \geq \frac{\Delta T}{2}$.

By Lemma 2 and Markov's inequality, w.p. ≥ 0.9 under $P_{\bar{m}}^{\text{ALG}'}$,

$$\begin{aligned}
 \frac{1}{T} \sum_{t=1}^{T\wedge \tau} g^M(a_t) &= \frac{1}{T} \sum_{t=1}^{T\wedge \tau} (r_*^M - r^M(a_t)) \\
 &\geq \frac{1}{T} \sum_{t=1}^{T\wedge \tau} (r_*^M - r_{\bar{m}}^M(a_t)) - \frac{1}{T} \sum_{t=1}^T |r^M(a_t) - r_{\bar{m}}^M(a_t)| \\
 \stackrel{\text{Lemma 2}}{\geq} &\frac{1}{T} \sum_{t=1}^{T\wedge \tau} g_{\bar{m}}^M(a_t) + \frac{T\wedge \tau}{T} (r_*^M - r_{\bar{m}}^M) - 10\varepsilon \\
 \stackrel{1)}{\geq} &\frac{1}{T} \sum_{t=1}^{T\wedge \tau} g_{\bar{m}}^M(a_t) + \frac{T\wedge \tau}{T} (0.99\Delta - \varepsilon) - 10\varepsilon,
 \end{aligned}$$

- If $\tau \geq T$, the RHS is $\geq 0.99\Delta - 11\varepsilon = \Omega(\Delta)$ for $\Delta > C\varepsilon$;
- If $\tau < T$, the RHS $\geq \frac{1}{T} \frac{\Delta T}{T_{\infty}} - 10\varepsilon = \Omega(\Delta)$ for $\Delta > C\varepsilon$.

$$\Rightarrow P_{\bar{m}}^{\text{ALG}'} \left(\frac{1}{T} \sum_{t=1}^{T\wedge \tau} g^M(a_t) = \Omega(\Delta) \right) \geq 0.9 \text{ in both cases.}$$

Since $TV(P_{\bar{m}}^{\text{ALG}'}, P_m^{\text{ALG}'}) \leq 0.1$ by Lemma 1, we get

$$P_m^{\text{ALG}'} \left(\frac{1}{T} \sum_{t=1}^{T\wedge \tau} g^M(a_t) = \Omega(\Delta) \right) \geq 0.8.$$

Finally, since ALG' and ALG coincide up to time $T\wedge \tau$, this gives the claimed result.

(In (Glasgow and Rakhlin '23), this stopping time argument establishes a stronger claim:

$$\inf_{\{a_1\}} \sup_M P_{\bar{m}}^{\text{ALG}'} \left(\frac{R(T)}{T} > ((1-\varepsilon) \text{dec}_{\varepsilon}(r) - C\varepsilon)_+ \right) = \Omega(1), \quad \varepsilon = \sqrt{\frac{c}{T}}$$

for any fixed $c_0 > 0$, where c, C depend on c_0 .)