# Lec 9: Advanced Fano's method

Yanjun Han

# Covering and packing

Let $(X, d)$ be a metric space, and $A \subseteq X$ be a compact set.

<u>Defn (covering)</u> $\{x_1, \cdots, x_n\} \subseteq X$ is an $\varepsilon$-covering (or $\varepsilon$-net) of $A$ if
$$A \subseteq \bigcup_{i=1}^{n} B(x_i; \varepsilon), \quad \text{with} \quad B(x; \varepsilon) = \{y \in X : d(x, y) \le \varepsilon\}.$$

<u>Defn (packing)</u> $\{a_1, \cdots, a_n\} \subseteq A$ is an $\varepsilon$-packing of $A$ if
$$\min_{i \ne j} d(a_i, a_j) > \varepsilon.$$

<u>Defn (covering and packing numbers)</u>
$$N(A, d, \varepsilon) = \min \{n : \exists \ \varepsilon\text{-covering of } A \text{ of size } n\}$$
$$M(A, d, \varepsilon) = \max \{m : \exists \ \varepsilon\text{-packing of } A \text{ of size } m\}$$

<u>Basic relationship.</u> $M(A, d, 2\varepsilon) \overset{①}{\le} N(A, d, \varepsilon) \overset{②}{\le} M(A, d, \varepsilon)$

(In other words, up to a multiplicative factor of 2 on $\varepsilon$, it's equivalent to consider covering or packing numbers.)

<u>Pf.</u> ① : If $M(A, d, 2\varepsilon) \ge N(A, d, \varepsilon) + 1$, by pigeonhole principle,
  $\exists$ two points $x, x'$ in a $(2\varepsilon)$-packing belong to the same ball $B(y; \varepsilon)$ in an $\varepsilon$-covering
  $\Rightarrow d(x, x') \le d(x, y) + d(x', y) \le 2\varepsilon$, a contradiction to $(2\varepsilon)$-packing.

② : If $a_1, \cdots, a_m$ is a maximal $\varepsilon$-packing of $A$, it must also be an $\varepsilon$-covering :
  if not, then $\exists \ a \in A$ s.t. $d(a, a_i) > \varepsilon \ \forall i \in [m]$
    $\Rightarrow \{a_1, \cdots, a_m, a\}$ is a larger $\varepsilon$-packing, a contradiction. $\boxed{6}$

# Bounding the covering/packing number.

1. <u>Volume bound</u> : let $\|\cdot\|$ be any norm on $\mathbb{R}^d$, $B = \{x : \|x\| \le 1\}$ be the unit ball.
   Then
$$\left(\frac{1}{\varepsilon}\right)^d \frac{\text{Vol}(A)}{\text{Vol}(B)} \overset{①}{\le} N(A, \|\cdot\|, \varepsilon) \le M(A, \|\cdot\|, \varepsilon) \overset{②}{\le} \left(\frac{2}{\varepsilon}\right)^d \frac{\text{Vol}(A + \frac{\varepsilon}{2}B)}{\text{Vol}(B)}.$$

**Pf.** ① : As $A \subseteq \overset{n}{\underset{i=1}{\cup}} B(x_i ; \varepsilon)$ for an $\varepsilon$-covering $\{x_1, \ldots, x_n\}$,

$$\text{Vol}(A) \leq \overset{n}{\underset{i=1}{\sum}} \text{Vol}(B(x_i ; \varepsilon)) = n \varepsilon^d \cdot \text{Vol}(B) \Rightarrow n \geq \left(\frac{1}{\varepsilon}\right)^d \frac{\text{Vol}(A)}{\text{Vol}(B)}.$$

② : As $\overset{m}{\underset{i=1}{\cup}} B(a_i ; \frac{\varepsilon}{2}) \subseteq A + \frac{\varepsilon}{2}B$, and the sets are disjoint under an $\varepsilon$-packing $\{a_1, \ldots, a_m\}$,

$$\text{Vol}(A + \frac{\varepsilon}{2}B) \geq \overset{m}{\underset{i=1}{\sum}} \text{Vol}(B(a_i ; \frac{\varepsilon}{2})) = m\left(\frac{\varepsilon}{2}\right)^d \text{Vol}(B) \Rightarrow m \leq \left(\frac{2}{\varepsilon}\right)^d \frac{\text{Vol}(A + \frac{\varepsilon}{2}B)}{\text{Vol}(B)}.$$

---

**Example 1.1.** If $A = \{x : \|x\| \leq 1\}$ is the unit ball under the same norm, then

$$\left(\frac{1}{\varepsilon}\right)^d \leq N(A, \|\cdot\|, \varepsilon) \leq \left(1 + \frac{2}{\varepsilon}\right)^d \leq \left(\frac{3}{\varepsilon}\right)^d \quad \text{for all} \quad 0 < \varepsilon \leq 1.$$

---

**Example 1.2 ( Gilbert-Varshamov bound).** If $A = \{0,1\}^d$ and $d_H(x, x') = \overset{d}{\underset{i=1}{\sum}} \mathbb{1}(x_i \neq x_i')$ is the Hamming distance, then for $1 \leq r \leq d-1$:

$$\frac{2^d}{\overset{r}{\underset{i=0}{\sum}} \binom{d}{i}} \leq M(A, d_H, r) \leq \frac{2^d}{\overset{r/2}{\underset{i=0}{\sum}} \binom{d}{i}}.$$

If $r = \rho d$ with $d \to \infty$, then by Stirling's approximation,

$$2^{d(1-h(\rho) + o(1))} \leq M(\{0,1\}^d, d_H, \rho d) \leq 2^{d(1-h(1/2) + o(1))}. \qquad \left( \begin{aligned} h(\rho) &= \rho \log_2 \frac{1}{\rho} \\ &+ (1-\rho)\log_2 \frac{1}{1-\rho} \end{aligned} \right)$$

**Pf.** Left inequality : $M(r) \geq N(r)$ and $\{0,1\}^d \subseteq \overset{n}{\underset{i=1}{\cup}} B(x_i ; r)$.

Right inequality : $\overset{m}{\underset{i=1}{\cup}} B(a_i ; \lfloor \frac{r}{2} \rfloor) \subseteq \{0,1\}^d$ and balls are disjoint.  ☒

---

**2. Sudakov minoration:** Let $w(A) = \mathbb{E} \underset{a \in A}{\sup} \langle a, Z \rangle$, $Z \sim N(\cdot, I_d)$ be the Gaussian width of $A$, then

$$w(A) \geq C \underset{\varepsilon > 0}{\sup} \ \varepsilon \sqrt{\log M(A, \|\cdot\|_2, \varepsilon)}.$$

( Two results needed to prove it:

① Slepian's lemma : let $X, Y$ be centered Gaussians in $\mathbb{R}^d$ with

$$\mathbb{E}[(Y_i - Y_j)^2] \leq \mathbb{E}[(X_i - X_j)^2] \qquad \forall i, j \in [d].$$

Then $\mathbb{E}[\max Y_i] \leq \mathbb{E}[\max X_i]$.

② Maximum of $n$ Gaussians : let $X_1, \ldots, X_n \overset{i.i.d.}{\sim} N(\cdot, 1)$, then

$$\mathbb{E}[\max X_i] = (1 + o(1)) \cdot \sqrt{2 \log n}. \qquad )$$

<u>Pf of Sudakov minoration</u>. Let $\{a_1, \dots, a_m\}$ be an optimal $\varepsilon$-packing of $A$.

$\quad$ Define $\quad X_i = \langle a_i, Z \rangle$ (with $Z \sim N(0, I_d)$), $\quad$ and $\quad Y_1, \dots, Y_m \overset{iid}{\sim} N(0, \frac{\varepsilon^2}{2})$.

$\quad$ Then $\quad \mathbb{E}[(Y_i - Y_j)^2] = \varepsilon^2 \leq \| a_i - a_j \|_2^2 = \mathbb{E}[(X_i - X_j)^2]$

$\quad\quad \implies w(A) \geq \mathbb{E}[\max_{i=1}^{m} X_i] \geq \mathbb{E}[\max_{i=1}^{m} Y_i] = \frac{\varepsilon}{\sqrt{2}} \cdot (1 + o(1)) \sqrt{2 \log m}$. $\qquad$ ∎


<u>Example 1.3</u>. When $A = B_1 = \{x : \|x\|_1 \leq 1\}$, then

$$w(A) = \mathbb{E} \sup_{\|x\|_1 \leq 1} \langle x, Z \rangle = \mathbb{E} \|Z\|_\infty \leq \sqrt{2 \log d}$$

$$\implies \log M(B_1, \|\cdot\|_2, \varepsilon) = O\left(\frac{\log d}{\varepsilon^2}\right).$$

In fact, it holds that

$$\log M(B_1, \|\cdot\|_2, \varepsilon) \asymp \begin{cases} d\left(1 + \log \frac{1}{\varepsilon^2 d}\right) & \text{if } \varepsilon \leq \frac{1}{\sqrt{d}} \quad \text{(volume bound is tight)} \\ \dfrac{1 + \log(\varepsilon^2 d)}{\varepsilon^2} & \text{if } \frac{1}{\sqrt{d}} < \varepsilon < 1 \quad \text{(Sudakov minoration nearly achieves} \\ & \qquad\qquad\qquad\qquad\quad \text{the tight upper bound)} \end{cases}$$


3. <u>Maurey's empirical method</u>: let $(H, \langle \cdot, \cdot \rangle)$ be an inner product space, and $T \subseteq H$ be a finite

$\quad$ set. Then $\qquad N(\text{conv}(T), \|\cdot\|, \varepsilon) \leq \binom{|T| + \lceil \frac{r^2}{\varepsilon^2} \rceil - 2}{\lceil \frac{r^2}{\varepsilon^2} \rceil - 1}, \quad 0 < \varepsilon \leq r,$

$\quad$ with $\qquad r = \inf_{y \in H} \sup_{x \in T} \|x - y\|$ is the <u>radius</u> of $T$.


<u>Pf</u>. We use a probabilistic argument. Let $T = \{t_1, \dots, t_m\}$, and $c \in H$ satisfy $r = \max_{i=1}^{m} \|t_i - c\|$.

$\quad$ Then for any $x \in \text{conv}(T)$, we have $x = \sum_{i=1}^{m} x_i t_i$ for some $x_i \geq 0$, $\sum_i x_i = 1$.

$\quad$ Let $Z$ be an $H$-valued RV s.t. $\mathbb{P}(Z = t_i) = x_i$ $1 \leq i \leq m$, then $x = \mathbb{E}[Z]$.

$\quad$ Let $Z_1, \dots, Z_n$ be i.i.d. copies of $Z$, and $\bar{Z} = \frac{1}{n+1}\left(c + \sum_{i=1}^{n} Z_i\right)$. Then

$$\mathbb{E}[\|\bar{Z} - x\|^2] = \frac{1}{(n+1)^2}\left(\underbrace{\|c - x\|^2}_{\substack{\leq r^2 \text{ by} \\ \text{convexity}}} + n\underbrace{\mathbb{E}[\|Z - x\|^2]}_{\substack{\leq \mathbb{E}[\|Z - c\|^2] \\ \leq r^2 \text{ as } \mathbb{E}Z = x}}\right) \leq \frac{r^2}{n+1}.$$

Consequently, if $n = \lceil \frac{r^2}{\varepsilon^2} \rceil - 1$, $\exists$ realization of $\bar{Z}$ s.t. $\|x - \bar{Z}\| \leq \varepsilon$. Meanwhile,

$$\bar{Z} \in \left\{ \frac{1}{n+1}\left(c + \sum_{i=1}^{m} n_i t_i\right) : \substack{n_i \geq 0 \\ \sum_i n_i = n} \right\} \text{ with cardinality } \binom{n + m - 1}{n}. \qquad \boxed{2}$$

<u>Example 1.3 cont'd</u>.   $B_1 = \text{conv}(\{\pm e_1, \cdots, \pm e_d\})$,  with radius 1.

By Maurey's empirical method,

$$\log N(B_1, \|\cdot\|_2, \varepsilon) \leq \log\left(\frac{2d + \lceil\frac{1}{\varepsilon^2}\rceil - 2}{\lceil\frac{1}{\varepsilon^2}\rceil - 1}\right) = O\left(\frac{1 + \log(\varepsilon^2 d)}{\varepsilon^2}\right) \text{ if } \frac{1}{\sqrt{d}} < \varepsilon < 1.$$

## 4. More results without proof.

① For $0 < p < q \leq \infty$, and $B_p := \{x \in \mathbb{R}^d : \|x\|_p \leq 1\}$, then

$$\log N(B_p, \|\cdot\|_q, \varepsilon) \asymp_{p,q} \begin{cases} \varepsilon^{-\frac{pq}{q-p}}(\log(d\varepsilon^{\frac{pq}{q-p}}) + 1) & \text{if } d^{\frac{1}{q} - \frac{1}{p}} \leq \varepsilon < 1. \\ d(\log\frac{1}{d\varepsilon^{\frac{pq}{q-p}}} + 1) & \text{if } \varepsilon < d^{\frac{1}{q} - \frac{1}{p}}. \end{cases}$$

② Let $N(A, B)$ be the smallest translations of $B$ that cover $A$. There exist universal constant $\alpha, \beta > 0$ such that, for any symmetric convex body $A$,

$$\frac{1}{\beta} \log N(B_2, \frac{\varepsilon}{\alpha} A^\circ) \leq \log N(A, \varepsilon B_2) \leq \beta \log N(B_2, \alpha\varepsilon A^\circ),$$

where $A^\circ = \{y : \sup_{x \in A} \langle x, y \rangle \leq 1\}$ is the polar body of $A$.

③ Let $H^s = \{f \in C^s([0,1]) : \|f^{(s)}\|_\infty \leq 1\}$, then

$$\log N(H^s, \|\cdot\|_p, \varepsilon) \asymp_p \varepsilon^{-\frac{1}{s}} \quad \text{for any } 1 \leq p \leq \infty.$$

④ Let $\mathcal{F}_M = \{f : [0,1] \to [0,1], f \text{ is non-decreasing}\}$, then

$$\log N(\mathcal{F}_M, \|\cdot\|_p, \varepsilon) \asymp_p \frac{1}{\varepsilon} \quad \text{for any } 1 \leq p < \infty.$$

⑤ Let $\mathcal{F}_c = \{f : [0,1] \to [0,1], f \text{ is convex}\}$, then

$$\log N(\mathcal{F}_c, \|\cdot\|_p, \varepsilon) \asymp_p \frac{1}{\sqrt{\varepsilon}} \quad \text{for any } 1 \leq p < \infty.$$

## Global Fano's method

Recall the steps of Fano's inequality:

① Find a pairwise separated set $\{\theta_1, \cdots, \theta_m\} \subseteq \Theta$ s.t. for all $i \neq j$,

$$\min_{a \in A} L(\theta_i, a) + L(\theta_j, a) \geq \Delta$$

② Try to upper bound $I(\theta; X)$ with $\theta \sim \text{Unif}(\{\theta_1, \cdots, \theta_m\})$ and $X | \theta \sim P_\theta$

③ If $I(\theta; X) < \frac{1}{2} \log m$, then the minimax risk satisfies $r^* = \Omega(\Delta)$.

Step ① is packing: if there is a metric $d(\theta, \theta')$ satisfies
$$\min_{a \in A} L(\theta, a) + L(\theta', a) \geq h(d(\theta, \theta')) \quad \text{for an increasing function } h: \mathbb{R}_+ \to \mathbb{R}_+,$$
then a $\delta$-packing $\{\theta_1, \ldots, \theta_M\}$ of $\Theta$ under $d$ satisfies the separation condition with $\Delta = h(\delta)$.

<span style="color:red">Q: Is there a general upper bound of $I(\theta; X)$, or more often, $I(\theta; X^n)$?</span>
<span style="color:red">A: This is possible using a covering of $(P_\theta)_{\theta \in \Theta}$ under KL!</span>

**Def.** For a family $P$ of distributions and $\varepsilon > 0$. let $N_{KL}(P, \varepsilon)$ be the smallest integer $n$
s.t. $\exists$ distributions $Q_1, \ldots, Q_n$ (not necessarily in $P$) satisfying
$$\sup_{P \in P} \min_{i \in [n]} D_{KL}(P \| Q_i) \leq \varepsilon^2.$$
<span style="color:red">($D_{KL}$ is <u>not</u> a metric; $Q_i$ in second argument)</span>

**Thm.** (Entropic upper bound of $I(\theta; X^n)$). Let $\theta \sim \pi$ with $\text{supp}(\pi) = \Theta_0$, and $X^n | \theta \sim P_\theta^{\otimes n}$.
Then
$$I(\theta; X^n) \leq \inf_{\varepsilon > 0} \left( n\varepsilon^2 + \log N_{KL}((P_\theta)_{\theta \in \Theta_0}, \varepsilon) \right).$$

**Pf.** Recall the "golden formula" in Lec 7:
$$I(\theta; X^n) = \min_{Q_{X^n}} \mathbb{E}_{\theta \sim \pi} \left[ D_{KL}(P_\theta^{\otimes n} \| Q_{X^n}) \right].$$
For an $\varepsilon$-covering of $(P_\theta)_{\theta \in \Theta_0}$, $Q_1, \ldots, Q_N$ with $N = N_{KL}((P_\theta)_{\theta \in \Theta_0}, \varepsilon)$, choose
$Q_{X^n} = \frac{1}{N} \sum_{i=1}^N Q_i^{\otimes n}$. Then for $\theta \sim \pi$,

$$D_{KL}\left( P_\theta^{\otimes n} \| \frac{1}{N} \sum_{i=1}^N Q_i^{\otimes n} \right) = \mathbb{E}_{P_\theta^{\otimes n}} \left[ \log \frac{P_\theta^{\otimes n}}{\frac{1}{N} \sum_{i=1}^N Q_i^{\otimes n}} \right]$$

$$\leq \mathbb{E}_{P_\theta^{\otimes n}} \left[ \min_{i \in [N]} \log \frac{P_\theta^{\otimes n}}{Q_i^{\otimes n}} + \log N \right] \quad \left( \textcolor{red}{\sum_i x_i \geq \max_i x_i} \right)$$

$$\leq \min_{i \in [N]} \mathbb{E}_{P_\theta^{\otimes n}} \left[ \log \frac{P_\theta^{\otimes n}}{Q_i^{\otimes n}} \right] + \log N$$

$$= \min_{i \in [N]} n D_{KL}(P_\theta \| Q_i) + \log N$$

$$\leq n\varepsilon^2 + \log N \quad \text{a.s. for } \theta \in \Theta_0. \qquad \boxed{\leftarrow}$$

Diagram of global Fano's method: for hyperparameters $\Theta_* \subseteq \Theta$, $\varepsilon$, $\delta > 0$:

① Find a metric $d(\theta, \theta')$ satisfying $\min_a L(\theta, a) + L(\theta', a) \geq h(d(\theta, \theta'))$ for an increasing non-negative function $h$, and find a $\delta$-packing of $\Theta_*$ under $d$.

② Find an $\varepsilon$-covering of $(P_\theta)_{\theta \in \Theta}$ under KL.

③ Apply Fano's method to conclude that

$$r^* \geq \frac{h(\delta)}{2}\left(1 - \frac{\log N_{KL}((P_\theta)_{\theta \in \Theta_*}, \varepsilon) + n\varepsilon^2 + \log 2}{\log M(\Theta_*, d, \delta)}\right).$$

Optimize over $(\Theta_*, \delta, \varepsilon)$ to make the lower bound as large as possible.

Example 2.1 (GLM). $X_1, \ldots, X_n \overset{i.i.d.}{\sim} N(\theta, I_d)$ with unknown $\theta \in \mathbb{R}^d$.

Target:
$$\inf_{\hat\theta} \sup_\theta \mathbb{E}_\theta[\|\hat\theta - \theta\|_p] \asymp_p \begin{cases} \frac{d^{1/p}}{\sqrt{n}} & 2 < p < \infty \\ \sqrt{\frac{\log d}{n}} & p = \infty \end{cases}.$$

Pf of lower bound. Choose $\Theta_* = \{\theta \in \mathbb{R}^d : \|\theta\|_2 \leq r\}$, then for any $\varepsilon, \delta > 0$,
global Fano gives

$$r^* \gtrsim \delta\left(1 - \frac{\log N_{KL}(\{N(\theta, I_d)\}_{\theta \in \Theta}, \varepsilon) + n\varepsilon^2 + \log 2}{\log M(\Theta_*, \|\cdot\|_p, \delta)}\right)$$

$$= \delta\left(1 - \frac{\log N(\Theta_*, \|\cdot\|_2, \sqrt{2}\varepsilon) + n\varepsilon^2 + \log 2}{\log M(\Theta_*, \|\cdot\|_p, \delta)}\right). \quad \left(\begin{array}{l} D_{KL}(N(\theta, I_d) \| N(\theta', I_d)) \\ = \frac{1}{2}\|\theta - \theta'\|_2^2 \end{array}\right)$$

Choice of $\varepsilon$: we choose $\varepsilon = \frac{r}{\sqrt{2}}$, so that $\log N = \log 1 = 0$.

Choice of $r/\delta$: for $p \in (2, \infty)$, choose $\frac{\delta}{r} = d^{\frac{1}{p} - \frac{1}{2}}$, so that
$$\log M(\Theta_*, \|\cdot\|_p, \delta) \asymp d.$$

For $p = \infty$, choose $\frac{\delta}{r} \asymp 1$, so that $\log M(\Theta, \|\cdot\|_\infty, \delta) \asymp \log d$.

Choice of $r$: now we have
$$r^* \gtrsim \begin{cases} rd^{\frac{1}{p} - \frac{1}{2}}\left(1 - \frac{C_1 nr^2 + \log 2}{C_2 d}\right) & \text{if } p \in (2, \infty), \\ r\left(1 - \frac{C_1 nr^2 + \log 2}{C_2 \log d}\right) & \text{if } p = \infty. \end{cases}$$

So $r = \begin{cases} \sqrt{\frac{d}{n}} & \text{for } p \in (2, \infty) \\ \sqrt{\frac{\log d}{n}} & \text{for } p = \infty \end{cases}$ gives $r^* \gtrsim \begin{cases} \frac{d^{1/p}}{\sqrt{n}} & 2 < p < \infty, \\ \sqrt{\frac{\log d}{n}} & p = \infty. \end{cases}$

**Example 2.2 (nonparametric density estimation)** $X_1, \cdots, X_n \overset{iind}{\sim} f$ on $[0,1]$ with $\|f^{(s)}\|_\infty \leq 1$.
(i.e. the function space is $H^s$)

Target:
$$\inf_{\hat{f}} \sup_{f \in F} \mathbb{E}_f \left[\|\hat{f}(X^n) - f\|_p\right] \gtrsim_p n^{-\frac{s}{2s+1}}, \quad p \in [1, \infty).$$

**Pf of lower bound.** Consider a subset $H_o^s \subseteq H^s$: $H_o^s = \{f \in H^s : f \geq \frac{1}{2} \text{ on } [0,1]\}$.

Then for $f, g \in H_o^s$, $D_{KL}(f\|g) \leq \chi^2(f\|g) \leq 2\|f-g\|_2^2$

$$\Rightarrow N_{KL}(H_o^s, \varepsilon) \leq N(H_o^s, \|\cdot\|_2, \frac{\varepsilon}{\sqrt{2}}) \leq N(H^s, \|\cdot\|_2, \frac{\varepsilon}{\sqrt{2}}).$$

By global Fano, for any $\varepsilon, \delta > 0$,

$$r^* \gtrsim \delta \left(1 - \frac{\log N_{KL}(H_o^s, \varepsilon) + n\varepsilon^2 + \log 2}{\log M(H_o^s, \|\cdot\|_p, \delta)}\right)$$

$$\geq \delta \left(1 - \frac{c_1 \varepsilon^{-\frac{1}{s}} + n\varepsilon^2 + \log 2}{c_2 \delta^{-\frac{1}{s}}}\right). \quad \text{(by metric entropy bounds for } H_o^s)$$

Choosing $\varepsilon \asymp \delta \asymp n^{-\frac{s}{2s+1}}$ gives $r^* = \Omega(n^{-\frac{s}{2s+1}})$. ☐

**Example 2.3 (Isotonic regression)** $X_1, \cdots, X_n \overset{i.i.d.}{\sim} P_X$, where (the known or unknown) $P_X$
has a bounded density on $[0,1]$. Conditioned on $X^n$, $Y_i \overset{ind}{\sim} N(f(X_i), 1)$ with
$$f \in F_M = \{f : [0,1] \to [0,1], f \text{ is increasing}\}.$$

Target:
$$\inf_{\hat{f}} \sup_{f \in F_M} \mathbb{E}_f \left[\|\hat{f} - f\|_p\right] \asymp_p n^{-\frac{1}{3}}, \quad \text{for all } p \in [1, \infty).$$

**Pf of lower bound.** Since $P_X$ has a bounded density,
$$D_{KL}(P_f \| P_{f'}) = \frac{1}{2}\|f - f'\|_{L^2(P_X)}^2 = O(1) \cdot \|f - f'\|_2^2$$
$$\Rightarrow N_{KL}((P_f)_{f \in F_M}, \varepsilon) \leq N(F_M, \|\cdot\|_2, \frac{\varepsilon}{O(1)}).$$

By global Fano:
$$r^* \gtrsim \delta \left(1 - \frac{\log N(F_M, \|\cdot\|_2, \frac{\varepsilon}{O(1)}) + n\varepsilon^2 + \log 2}{\log M(F_M, \|\cdot\|_p, \delta)}\right)$$

$$\geq \delta \left(1 - \frac{\frac{c_1}{\varepsilon} + n\varepsilon^2 + \log 2}{1/\delta}\right). \quad (\log N(F_M, \|\cdot\|_p, \varepsilon) \asymp_p \frac{1}{\varepsilon})$$

Choosing $\varepsilon \asymp n^{-1/3}$ and $\delta \asymp n^{-1/3}$, we obtain $r^* = \Omega(n^{-1/3})$. ☐

Example 2.4 (Convex regression) Same setting as Example 2.3, but with $\mathcal{F}_m$ replaced by
$$\mathcal{F}_c = \{f : [0,1] \to [0,1]. \; f \text{ is convex}\}.$$

Target:
$$\inf_{\hat{f}} \sup_{f \in \mathcal{F}_c} \mathbb{E}_f\left[\|\hat{f} - f\|_p\right] \asymp_p n^{-\frac{2}{5}}, \qquad p \in [1, \infty).$$

__Pf of lower bound__, similar to Example 2.3, now with $\log N(\mathcal{F}_c, \|\cdot\|_p, \varepsilon) \asymp \frac{1}{\sqrt{\varepsilon}}$. ☐

Example 2.5 (Sparse linear regression) $y \sim N(X\theta, I_n)$ with fixed design $X \in \mathbb{R}^{n \times d}$,
where all singular values of $X$ are $O(\sqrt{n})$. The unknown parameter $\theta \in \mathbb{R}^d$ is sparse:
$$\|\theta\|_q \leq R \qquad \text{for some } q \in (0,1).$$

Target:
$$\inf_{\hat{\theta}} \sup_{\|\theta\|_q \leq R} \mathbb{E}_\theta\left[\|\hat{\theta} - \theta\|_p\right] \asymp_{p,q} R^{\frac{q}{p}}\left(\frac{\log d}{n}\right)^{\frac{p-q}{2p}} \quad \text{for small enough}$$
$$R < f(n,d).$$

__Pf of lower bound__.

1. $L_p$-packing of $B_q(R) = \{\theta \in \mathbb{R}^d : \|\theta\|_q \leq R\}$:
$$\log M(B_q(R), \|\cdot\|_p, \delta) \asymp \left(\frac{R}{\delta}\right)^{\frac{pq}{p-q}} \log d \qquad \text{if } \delta \gg R d^{\frac{1}{p} - \frac{1}{q}}.$$

2. KL covering of $\mathcal{P} = \{N(X\theta, I_n) : \|\theta\|_q \leq R\}$:
$$D_{KL}\left(N(X\theta, I_n) \| N(X\theta', I_n)\right) = \frac{1}{2}\|X(\theta - \theta')\|_2^2 = O(n) \cdot \|\theta - \theta'\|_2^2$$
$$\implies \log N_{KL}(\mathcal{P}, \varepsilon) \leq \log N\left(B_q(R), \|\cdot\|_2, \frac{\varepsilon}{O(\sqrt{n})}\right)$$
$$\asymp \left(\frac{\sqrt{n}R}{\varepsilon}\right)^{\frac{2q}{2-q}} \log d \qquad \text{if } \varepsilon \gg R\sqrt{n}\, d^{\frac{1}{2} - \frac{1}{q}}.$$

Now choosing
$$\varepsilon \asymp n^{\frac{1}{4}} R^{\frac{q}{2}} (\log d)^{\frac{2-q}{4}}$$
$$\delta \asymp R^{\frac{q}{p}}\left(\frac{\log d}{n}\right)^{\frac{p-q}{2p}},$$
then
$$\log M(\delta) \asymp R^q n^{\frac{q}{2}} (\log d)^{1 - \frac{q}{2}}.$$
$$\log N_{KL}(\varepsilon) \asymp \varepsilon^2 \asymp R^q n^{\frac{q}{2}} (\log d)^{1 - \frac{q}{2}}.$$
and global Fano gives the result. ☐

<u>Special topic</u>: generalized Fano with $\chi^2$-informativity.

Since the proof of Fano is simply DPI, replacing KL by other $f$-divergences also leads to meaningful Bayes risk lower bounds.
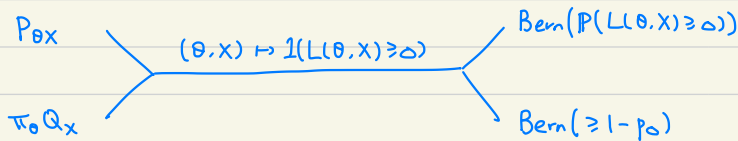
<u>Thm</u>. For $\theta \sim \pi$, it holds that
$$\mathbb{P}(L(\theta, X) \geq \Delta) \geq 1 - p_\Delta - \sqrt{p_\Delta \cdot I_{\chi^2}(\theta ; X)},$$
where $p_\Delta = \sup_{\hat{a}} \pi(L(\theta, a) < \Delta)$ is the small-ball probability, and
$$I_{\chi^2}(\theta ; X) = \inf_{Q_X} \chi^2(P_{\theta X} \| \pi_\theta Q_X) = \inf_{Q_X} \mathbb{E}_{\theta \sim \pi}[\chi^2(P_{X|\theta} \| Q_X)]$$
is the <u>$\chi^2$-informativity</u>.

<u>Pf</u>. Apply DPI to:

$$\begin{array}{ccc} P_{\theta X} & & \text{Bern}(\mathbb{P}(L(\theta, X) \geq \Delta)) \\ & \searrow \quad (\theta, X) \mapsto \mathbb{1}(L(\theta, X) \geq \Delta) \quad \nearrow & \\ & & \\ \pi_\theta Q_X & \nearrow \qquad\qquad\qquad \searrow & \text{Bern}(\geq 1 - p_\Delta) \end{array}$$

We get:
$$\chi^2(P_{\theta X} \| \pi_\theta Q_X) \geq \chi^2(\text{Bern}(\mathbb{P}(L(\theta, X) \geq \Delta)) \| \text{Bern}(\geq 1 - p_\Delta))$$
$$\geq \frac{(\mathbb{P}(L(\theta, X) \geq \Delta) - (1 - p_\Delta))^2}{p_\Delta(1 - p_\Delta)} \quad \text{if } \mathbb{P}(L(\theta, X) \geq \Delta) \leq 1 - p_\Delta.$$

Taking inf over $Q_X$ and rearranging gives the result. $\qquad\qquad\qquad$ ③

Similarly, we have an entropic upper bound of $I_{\chi^2}(\theta ; X)$ based on $\chi^2$-covering.

<u>Thm</u>. Let $\mathcal{P} = (P_\theta)_{\theta \in \Theta}$ and $\text{supp}(\pi) \subseteq \Theta$. Then for $\theta \sim \pi$,
$$I_{\chi^2}(\theta ; X) + 1 \leq \inf_{\varepsilon > 0} (1 + \varepsilon^2) N_{\chi^2}(\mathcal{P}, \varepsilon),$$
where $N_{\chi^2}(\mathcal{P}, \varepsilon) = \min \left\{ n : \min_{Q_1, \cdots, Q_n} \sup_{P \in \mathcal{P}} \min_{i \in [n]} \chi^2(P \| Q_i) \leq \varepsilon^2 \right\}.$

<u>Pf</u>. Exercise.

Example 3.1 (Gaussian model with uniform prior). Let $X \sim N(\theta, I_d)$ with $\theta \sim \text{Unif}(B_2(R))$.

Target: $\quad r_\pi^* = \inf_{\hat\theta} \mathbb{E}_{\theta \sim \pi}[\|\hat\theta - \theta\|_2^2] \asymp d \quad$ if $\quad R = \Omega(\sqrt{d})$.

**Failure of mutual information.** For $\Delta \in (0, R)$, the small ball prob. is

$$p_\Delta = \sup_a \pi\left(\|\theta - a\|_2^2 \le \Delta^2\right) = \left(\frac{\Delta}{R}\right)^d.$$

For the mutual info, the entropic upper bound gives

$$I(\theta; X) \le \inf_{\varepsilon > 0}\left(\log N(B_2(R), \|\cdot\|_2, \varepsilon) + \varepsilon^2\right)$$
$$\le \inf_{\varepsilon > 0}\left(d\log \frac{3R}{2} + \varepsilon^2\right) \sim d\log\frac{R}{\sqrt{d}} \quad \text{if} \quad R \gg \sqrt{d}$$

Therefore, Fano gives that

$$r_\pi^* \gtrsim \sup_{\Delta > 0} \Delta^2 \cdot \left(1 - \frac{d\log\frac{R}{\sqrt{d}} + \log 2}{d\log\frac{R}{\Delta}}\right).$$

Usually we make

$$d\log\frac{R}{\sqrt{d}} = (1-\rho) d\log\frac{R}{\Delta} \implies \Delta = d^{\frac{1}{2(1-\rho)}} R^{-\frac{\rho}{1-\rho}}, \text{ for some constant } \rho > 0$$
$$\implies r_\pi^* = \Omega\left(d^{\frac{1}{1-\rho}} R^{-\frac{2\rho}{1-\rho}}\right) = \Omega\left(d \cdot \left(\frac{d}{R^2}\right)^{\frac{\rho}{1-\rho}}\right) \text{ is weaker than } \Omega(d).$$

**Pf using $\chi^2$-info.** The entropic upper bound gives

$$I_{\chi^2}(\theta; X) + 1 \le \inf_{\varepsilon > 0}(1 + \varepsilon^2) N\left(B_2(R), \|\cdot\|_2, \sqrt{\log(1+\varepsilon^2)}\right) \quad \left(\chi^2(N(\theta, I) \| N(\theta', I))\right)$$
$$= e^{\|\theta - \theta'\|_2^2} - 1)$$
$$\le \inf_{\varepsilon > 0}(1 + \varepsilon^2)\left(\frac{3R}{\sqrt{\log(1+\varepsilon^2)}}\right)^d = \exp\left(d\log\frac{O(1) \cdot R}{\sqrt{d}}\right) \quad \text{for } R > C\sqrt{d}.$$
$$\text{by choosing } 1 + \varepsilon^2 = e^d.$$

Therefore, generalised Fano gives

$$r_\pi^* \gtrsim \sup_{\Delta > 0} \Delta^2\left(1 - \left(\frac{\Delta}{R}\right)^d - \sqrt{\underbrace{\left(\frac{\Delta}{R}\right)^d \cdot \exp\left(d\log\frac{O(1) \cdot R}{\sqrt{d}}\right)}_{\le \frac{1}{2} \text{ for } \Delta = c\sqrt{d} \text{ with a small constant } c}}\right)$$

$$= \Omega(d).$$

Example 3.2 (ridge bandits). $r_t \sim N(f(\langle \theta^*, a_t \rangle), 1)$ for $\theta^* \sim \text{Unif}(S^{d-1})$,

$f$: known link function $[-1,1] \to \mathbb{R}$, increasing, and $f(0) = 0$.

Target: Define a recursive sequence with a large constant $C > 0$:

$$\varepsilon_1 = C \sqrt{\frac{\log(1/\delta)}{d}} \quad, \quad \varepsilon_{t+1}^2 = \varepsilon_t^2 + \frac{C}{d} g(\varepsilon_t)^2 \quad (g(x) := \max\{|f(x)|, |f(-x)|\})$$

Then for any interactive learner,

$$\mathbb{P}\left( |\langle \theta^*, a_s \rangle| \leq \varepsilon_s \text{ for all } 1 \leq s \leq t \right) \geq 1 - t\delta.$$

Remark: ① The sequence $\{\varepsilon_t\}$ is a <u>pointwise</u> upper bound on the learning trajectory of <u>any</u> algorithm

② The growth $\varepsilon_{t+1}^2 - \varepsilon_t^2$ increases with $t$: interactive learning becomes faster and faster!

<u>Intuition</u>. Let $I_t = I(H_t; \theta^*) := I(a^t, r^t; \theta^*)$. Then

$$I_{t+1} - I_t = I(\theta^*; r_{t+1} \mid H_t, a_{t+1})$$

$$\leq \mathbb{E}\left[ D_{KL}\left( N(f(\langle \theta^*, a_{t+1} \rangle), 1) \,\|\, N(0,1) \right) \right] \quad (\text{Golden formula})$$

$$= \frac{1}{2} \mathbb{E}\left[ f(\langle \theta^*, a_{t+1} \rangle)^2 \right].$$

We aim to upper bound this information gain. A key observation is that,

$$I(\theta^*; a_{t+1}) \leq I(\theta^*; H_t) = I_t,$$

so $a_{t+1}$ is "constrained" in information, and we expect $\langle \theta^*, a_{t+1} \rangle$ to be small. The intuition is that:

$$I(\theta^*; a) \leq d\varepsilon^2 \implies |\langle \theta^*, a \rangle| \leq \varepsilon \text{ w.h.p.} \qquad (*)$$

If $(*)$ were true, we'll get the recursion by the correspondence $I_t \asymp d\varepsilon_t^2$. However, mutual info is not strong enough to ensure $(*)$: Fano only gives

$$\mathbb{P}\left( |\langle \theta^*, a \rangle| \leq \varepsilon \right) \geq 1 - \underline{\frac{I(\theta^*; a) + \log 2}{c\, d\varepsilon^2}}$$

not small enough to apply union bound!

<u>Pf using $\chi^2$-info</u>. Let $E_t = \bigcap_{s \leq t} \{|\langle \theta^*, a_s \rangle| \leq \varepsilon_s\}$. Define a slight variant of $\chi^2$-info as

$$I_{\chi^2}(X; Y \mid E) = \inf_{Q_Y} \chi^2(P_{XY|E} \| P_X Q_Y).$$

then we can still get

$$P(|\langle \theta^*, a \rangle| \leq \varepsilon \mid E) \geq 1 - c_1 \underbrace{e^{-c_0 d \varepsilon^2}}_{\text{for fixed } a,} \sqrt{I_{\chi^2}(\theta^*; a \mid E) + 1}.$$

$$\color{red}\text{for fixed } a,$$
$$\color{red}P(|\langle \theta^*, a \rangle| \leq \varepsilon) \leq e^{-c_0 d \varepsilon^2}.$$

The crux of the proof is to establish the following recursion:

$$I_{\chi^2}(\theta^*; H_t \mid E_T) + 1 \leq \frac{e^{g(\varepsilon_t)^2}}{P(E_t \mid E_{T-1})^2}\left(I_{\chi^2}(\theta^*; H_{t-1} \mid E_{t-1}) + 1\right). \qquad (*)$$

If $(*)$ holds, then $I_{\chi^2}(\theta^*; H_t \mid E_T) \leq \prod_{s=1}^{t} \frac{e^{g(\varepsilon_s)^2}}{P(E_s \mid E_{s-1})^2} = \frac{1}{P(E_t)^2} \exp\left(\sum_{s \leq t} g(\varepsilon_s)^2\right)$, so

$$P(E_{t+1} \mid E_T) \geq 1 - c_1 e^{-c_0 d \varepsilon_{t+1}^2} \sqrt{I_{\chi^2}(\theta^*; a_{t+1} \mid E_T) + 1}$$

$$\geq 1 - c_1 e^{-c_0 d \varepsilon_{t+1}^2} \sqrt{I_{\chi^2}(\theta^*; H_t \mid E_T) + 1} \qquad \color{red}(DPI)$$

$$\geq 1 - \frac{c_1}{P(E_t)} \exp\left(\underbrace{-c_0 d \varepsilon_{t+1}^2 + \frac{1}{2} \sum_{s \leq t} g(\varepsilon_s)^2}_{\color{red}\text{recursion ensures that } \leq -c_0 d \varepsilon_t^2 \leq -c' \log(\frac{1}{\delta})}\right)$$

$$\implies P(E_{t+1}) = P(E_T) \cdot P(E_{t+1} \mid E_T) \geq P(E_T) - \delta. \qquad \boxed{\triangle}$$

<u>Pf of $(*)$</u>.

$$I_{\chi^2}(\theta^*; H_t \mid E_T) + 1 = \inf_{Q_{H_t}} \int \frac{P(\theta^*, H_t \mid E_T)^2}{\pi(\theta^*) Q_{H_t}(H_t)} d\theta^* da^t dr^t$$

$$\leq \inf_{Q_{H_{t-1}}} \int \frac{\left[\frac{1(E_t)}{P(E_T)} \pi(\theta^*) \prod_{s=1}^{t}\left(P_s(a_s \mid H_{s-1}) \varphi(r_s - f(\langle \theta^*, a_s \rangle))\right)\right]^2}{\pi(\theta^*) Q_{H_{t-1}}(H_{t-1}) \cdot P_t(a_t \mid H_{t-1}) \varphi(r_t)} d\theta^* da^t dr^t \qquad \color{red}\leq g(\varepsilon_t)^2 \text{ on } E_T$$

$$= \inf_{Q_{H_{t-1}}} \int \frac{\left[\underbrace{\frac{1(E_t)}{P(E_T)}}_{\color{red}\leq \frac{1(E_{t-1})}{P(E_{t-1})} \cdot \frac{1}{P(E_T \mid E_{t-1})}} \pi(\theta^*) \prod_{s=1}^{t-1}\left(P_s(a_s \mid H_{s-1}) \varphi(r_s - f(\langle \theta^*, a_s \rangle))\right)\right]^2}{\pi(\theta^*) Q_{H_{t-1}}(H_{t-1})} \cdot P_t(a_t \mid H_{t-1}) e^{\color{red}f(\langle \theta^*, a_s \rangle)^2} d\theta^* da^t dr^{t-1}$$

$$\leq \frac{\exp(g(\varepsilon_t)^2)}{P(E_t \mid E_{t-1})^2}\left(I_{\chi^2}(\theta^*; H_{t-1} \mid E_{t-1}) + 1\right) \qquad \boxed{\triangle}$$