

# Minimax risk theory

# 1. Problem Setting

Data:  $X_1, X_2, \dots, X_n \sim P, P \in \mathcal{P}$  family

☆  $\hat{\theta}(X_1, X_2, \dots, X_n) \xrightarrow{\text{estimate}} \theta(P)$

Wondering:

☆  $R_n = \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_{X_i \sim P} [d(\hat{\theta}, \theta(P))]$

"minimax risk"

(For lower regret bound, we should firstly find a  $d(\cdot, \cdot)$  that lower bounds Reg.)

$d$  satisfies the "triangular inequality".

Example 1:  $\mathcal{P} = \{N(\theta, 1), \theta \in \mathbb{R}\}$

want to know  $R_n = \inf_{\hat{\theta}} \sup_{\theta} \mathbb{E}[(\hat{\theta} - \theta)^2]$

△ why?  $\Rightarrow \hat{\theta} = \frac{1}{n} \sum_{i=1}^n X_i$

Note: What we really care about is the problem hardness than the estimator.

minimax risk is useful to represent the hardness even for non-parametric family where we don't even

have a close-form solution.

To ~~know~~  $R_n$  is not necessary.

↓  
bound

as  $L_n \leq R_n \leq U_n$  is necessary,

(best result:  $U_n = O(L_n)$ )

Define 2 metrics:

$\Delta$  is  $KL(P, Q) = \int p \log \frac{p}{q}$  or  $= \sum_i p_i \log \frac{p_i}{q_i}$

Property 1:  $KL(P^n, Q^n) = n \cdot KL(P, Q)$  (for iid)

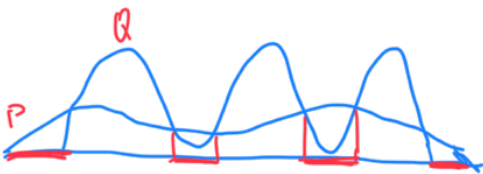
Property 1':  $KL(P^T, Q^T) = \sum_{t=1}^T KL(P[X_t | X^{t-1}], Q[X_t | X^t])$   
called "Chain rule of Relative entropy" (for martingale)

it's  $TV(P, Q) = \sup_{A \text{ event}} |P(A) - Q(A)| = \frac{1}{2} \int |P - Q|$

Homework

Proof:  $\int |P - Q| = \int_{P(x) \geq Q(x)} (P(x) - Q(x)) dx$   $C$

$+ \int_{P(x) < Q(x)} (Q(x) - P(x)) dx$   $C'$



Notice that  $C - C' = \int (P(x) - Q(x)) dx$

$$\begin{aligned}
 & + \int_C (P(x) - Q(x)) dx \\
 & = \int (P(x) - Q(x)) dx \\
 & = 1 - 1 \\
 & = 0
 \end{aligned}$$

$$\Rightarrow C = C'$$

$$\text{Also, } \sup_A |P(A) - Q(A)| = \max \langle C, C' \rangle$$

$$\Rightarrow \int |P - Q| = 2 \sup_A |P(A) - Q(A)|$$

↑ since  $C$  is the sum of all positive items, and  $C'$  is the sum of all negative items.

$$\text{Property 2: } TV(P, Q) \leq \sqrt{KL(P, Q)}$$

$$\text{Recall: } R_n = \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_P [d(\hat{\theta}, \theta(P))]$$

hard to calculate

△ Reduction ①: "Finite Covering"

discretize distribution family  $\mathcal{P}$

Pick a finite set  $M = \{P_1, P_2, \dots, P_N\} \subset \mathcal{P}$

Then:

$$R_n = \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_P [d(\hat{\theta}, \theta(P))]$$

- A

$$\geq \inf_{\hat{\theta}} \sup_{P_j \in \mathcal{M}} \mathbb{E}_{P_j} [d(\hat{\theta}, \theta_j)], \theta_j = \theta(P_j)$$

Remark: if  $\mathcal{M}$  can be a good representative of family  $\mathcal{P}$ , then this lower bound would approach  $R_n$ . Both "Representable" and "Indistinguishable"

Define:  $S = \min_{j \neq k} d(\theta_j, \theta_k)$ , then

Reduction (2): by Markov

$$P_j [d(\hat{\theta}, \theta_j) > \frac{S}{2}] \leq \frac{2}{S} \cdot \mathbb{E}_{P_j} [d(\hat{\theta}, \theta_j)], \forall j \in [N]$$

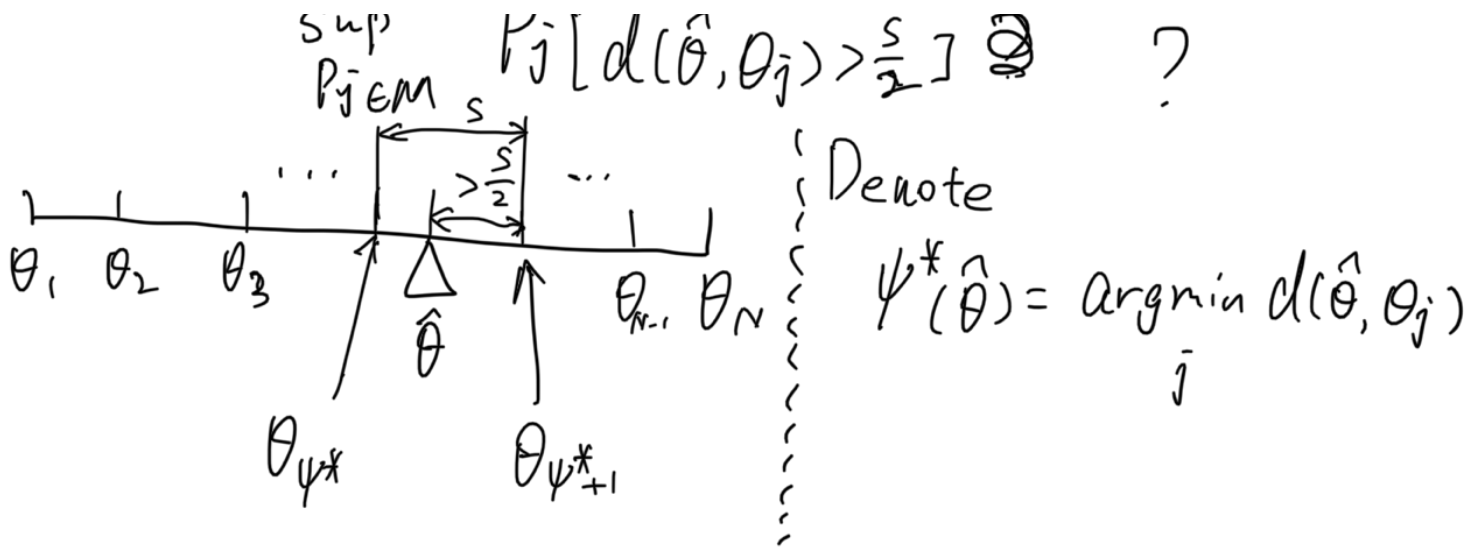
$$\Rightarrow \mathbb{E}_{P_j} [d(\hat{\theta}, \theta_j)] \geq \frac{S}{2} \cdot P_j [d(\hat{\theta}, \theta_j) > \frac{S}{2}] \quad \text{"1, 2, \dots, N"}$$

$$\Rightarrow R_n \geq \inf_{\hat{\theta}} \sup_{P_j \in \mathcal{M}} \mathbb{E}_{P_j} [d(\hat{\theta}, \theta_j)]$$

$$\geq \frac{S}{2} \cdot \inf_{\hat{\theta}} \sup_{P_j \in \mathcal{M}} P_j [d(\hat{\theta}, \theta_j) > \frac{S}{2}]$$

Still hard to compute, we need further lower bound

How to lower bound



Then: for any  $j \neq \psi^*$ , we have

$$d(\hat{\theta}, \theta_j) \geq \frac{s}{2}$$

Since

$$\begin{aligned}
 s &\leq d(\theta_j, \theta_{\psi^*}) \\
 &\leq d(\theta_j, \hat{\theta}) + d(\hat{\theta}, \theta_{\psi^*}) \\
 &\leq d(\theta_j, \hat{\theta}) + d(\hat{\theta}, \theta_j) \\
 &= 2d(\hat{\theta}, \theta_j)
 \end{aligned}$$

$\Rightarrow$  Reduction (3):

$$\begin{aligned}
 P_j [d(\hat{\theta}, \theta_j) \geq \frac{s}{2}] &\geq P_j [\psi^*(\hat{\theta}) \neq j] \\
 &\geq \inf_{\psi(\hat{\theta})} P_j [\psi(\hat{\theta}) \neq j] \\
 &\geq \inf_{\psi(x_1, x_2, \dots, x_n)} P_j [\psi \neq j], \forall \psi: X^n \rightarrow [N]
 \end{aligned}$$

Remark: Notice that  $\psi^*(\hat{\theta})$  is sth you choose, related to  $\hat{\theta}$  but  $j$  is chosen by the adversary.

$\hat{\theta}$  is chosen by the sup, or, say, by the nature

$$\Rightarrow R_n \geq \frac{s}{2} \cdot \inf_{\hat{\theta}} \sup_{P_j \in \mathcal{M}} P_j [d(\hat{\theta}, \theta_j) > \frac{s}{2}]$$

★  $\geq \frac{s}{2} \cdot \inf_{\psi} \sup_{P_j \in \mathcal{M}} P_j [\psi \neq j]$

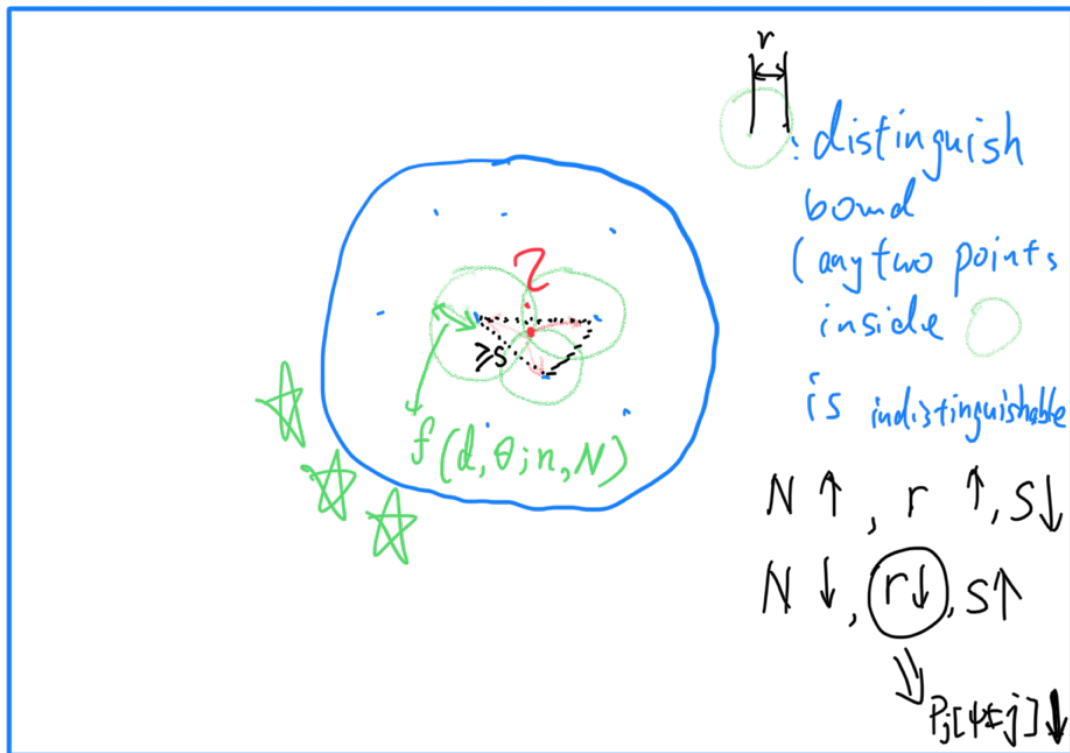
Remember:

$$s = \min_{j \neq k} d(\theta_j, \theta_k)$$

Compete

a minimax game

Illustration:

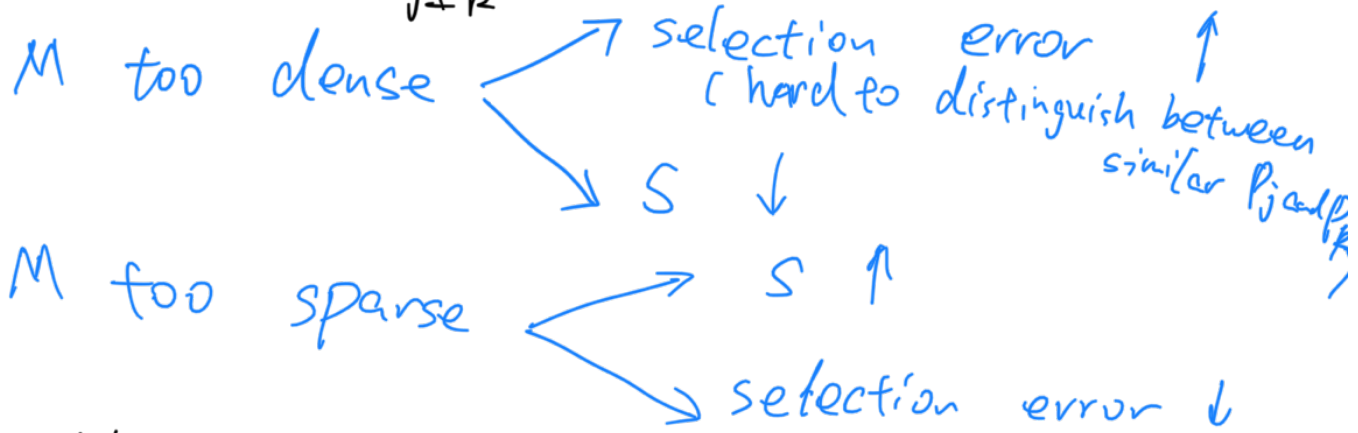


In the following, we focus on:

★  $R_n = \inf \sup_{P_j \in \mathcal{M}} \mathbb{E} [d(\hat{\theta}, \theta_j)] > \frac{s}{2} \inf \max_{P_j} P_j [i \neq j]$

$$\hat{\theta} = \arg \min_{\theta \in \mathcal{P}} \sum_{j \in M} P_j \psi(\theta, \theta_j)$$

Recall that  $S = \min_{j \neq k} d(\theta_j, \theta_k)$  is related to  $M$



$\rightarrow$  How to characterize the "distinguishability"?

i.e., how to calculate  $\inf_{\psi} \max_j P_j [\psi \neq j]$

3 Standard methods:

- Le Cam
- Fan
- Tsybakov

1: Le Cam: distinguish between 2 guys:

$$M = \{P_0, P_1\}$$

$\delta$  Theorem 4 (Le Cam): for any  $P_0, P_1 \in \mathcal{P}_{\text{prob}}$

$$R_n = \inf \sup \mathbb{E}[d(\hat{\theta}, \theta(P))] \geq \frac{S}{8} \cdot e^{-n \cdot KL(P_0, P_1)}$$



$\theta \in \mathcal{P}$

(Here  $s = d(\theta \in \mathcal{P}_0, \theta \in \mathcal{P}_1)$ )

★ Corollary : if  $KL(P_0, P_1) \leq \log \frac{2}{n}$ , then

$$R_n = \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E} [d(\hat{\theta}, \theta(P))] \geq \frac{s}{16}$$

Notice that LHS is nothing to do with  $P_0$  and  $P_1$ ,  
 Therefore, just find 2 distributions in  $\mathcal{P}$  that  
 is  $\log \frac{2}{n}$  apart, and then they are hard to  
 distinguish, and this in return ensures the  
 lower bound <sup>inequality</sup> to be true.

Proof of Le Cam's Theorem: (maybe show later?)

$$R_n \geq \frac{s}{2} \cdot \inf_{\psi} \max \{ P_0(\psi=1), P_1(\psi=0) \}$$

here  $\psi: X \rightarrow \{0, 1\}$

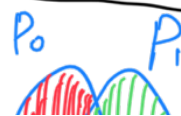
$$\geq \frac{s}{2} \inf_{\psi}$$

$$\frac{P_0(\psi=1) + P_1(\psi=0)}{2}$$

Neyman-Pearson's Test

$$\psi^*(x) = \begin{cases} 0 & \text{if } P_0(x) \geq P_1(x) \\ 1 & \text{if } P_1(x) > P_0(x) \end{cases}$$

"Type 1 error"



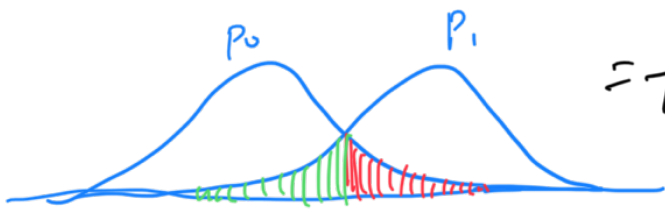
Lemma: for any  $\psi$ , we have

$$P_0(\psi=1) + P_1(\psi=0) \geq P_0(\psi^*=1) + P_1(\psi^*=0)$$

[Optimality of Neyman-Pearson's Test]

(No proof)

$$\geq \frac{S}{2} \cdot \frac{1}{2} (P_0(\psi^*=1) + P_1(\psi^*=0))$$



$$= \frac{S}{4} \cdot \int_{P_1(x) \geq P_0(x)} P_0(x) dx + \int_{P_1(x) < P_0(x)} P_1(x) dx$$

$$= \frac{S}{4} \int \min\{P_0(x), P_1(x)\} dx$$

Le Cam's Lemma:

$$\inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_P[d(\hat{\theta}, \theta(P))] \geq \frac{S}{4} \cdot \int \min\{P_0(x), P_1(x)\} dx$$

Lemma 7: for any  $P$  and  $Q$ , we have:

$$\int \min\{P(x), Q(x)\} dx \geq \frac{1}{2} e^{-KL(P, Q)}$$

(No proof) (Hint: by Jensen's Inequality)

With 1. and 7 we prove Thm 4

Example of Le Cam:

$$Y = m(X_i) + \varepsilon, \quad X_i \sim \text{Unif}(0, 1), \quad \varepsilon_i \sim N(0, 1)$$

$$\text{Let } m(\cdot) \in \mathcal{M} = \{m: |m(y) - m(x)| \leq L|x - y|, \quad x, y \in [0, 1]\}$$

(L-lipschitz family)

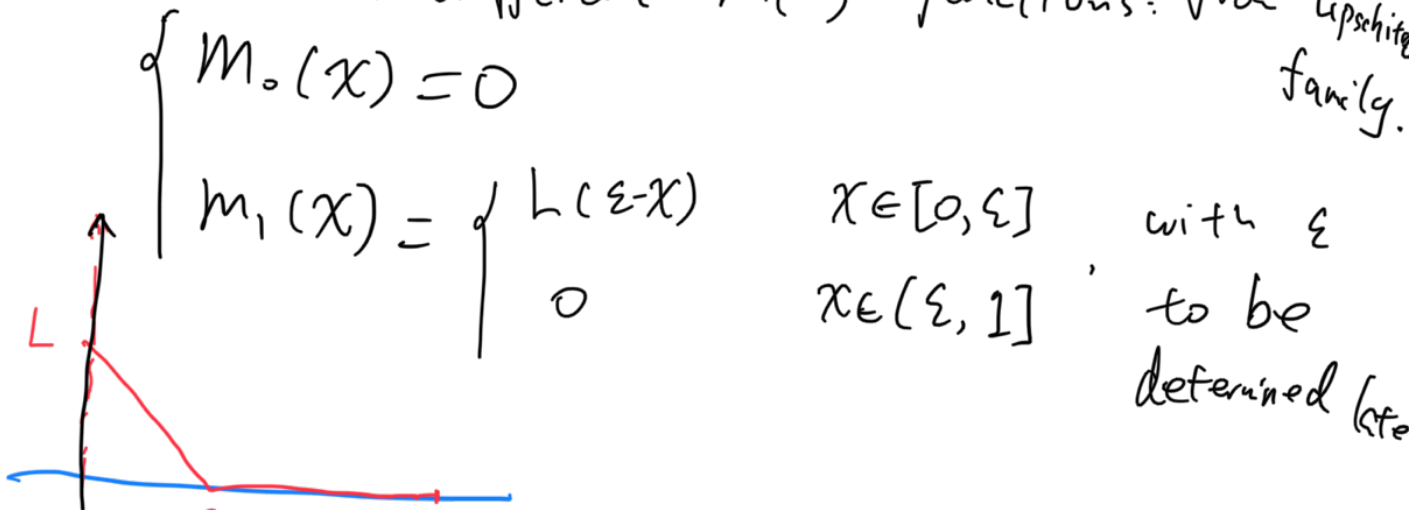
We want to bound the minimax risk of

$$\theta := m(0)$$

loss function:  $d(\theta_0, \theta_1) = |\theta_0 - \theta_1|$

Distributions:  $P(x, y) = \underbrace{\phi(x)}_{\text{uniform}} \underbrace{P(y|x)}_{\sim N(m(x), 1)}$

Now, in order to apply Le Cam's method, we have to choose 2 different  $m(\cdot)$  functions: from Lipschitz family.



0  $\epsilon$  1  $m_0$

$$\Rightarrow KL(P_0, P_1) = \iint P_0(x, y) \log \frac{P_0(x, y)}{P_1(x, y)} dx dy$$
$$= \int_0^\epsilon KL(N(0, 1), N(m(x), 1)) dx$$

Note:

$$KL(N(\mu_0, \sigma^2), N(\mu_1, \sigma^2)) = \frac{(\mu_1 - \mu_0)^2}{2\sigma^2}$$

$$= \int_0^\epsilon \frac{L^2(\epsilon - x)^2}{2} dx$$

$$= \frac{L^2 \epsilon^3}{6} \leftarrow \frac{\log 2}{n}$$

Let  $\epsilon = \left(\frac{c}{n}\right)^{\frac{1}{3}}$ , and we have:

$|m_1(0) - m_2(0)| \geq \left(\frac{c}{n}\right)^{\frac{1}{3}}$  according to Le Cam's Lemma, i.e.,

$\inf_{\hat{m}(0)} \sup_{m \in \mathcal{M}} \mathbb{E}[|\hat{m}(0) - m(0)|] \geq \left(\frac{c}{n}\right)^{\frac{1}{3}}$ , which is a lower error bound.

Note: This is also an upper bound.  $\Rightarrow$  minimax risk.

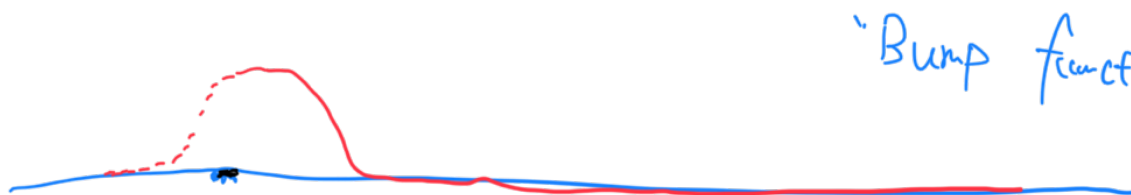
Note.2: for high-dimensional case:  $\dots = \Omega\left(\left(\frac{1}{n}\right)^{\frac{1}{d+2}}\right)$

---

Recall: how we constructed  $m_n(x)$ ?



for Lipschitz



"Bump function" for high-order

-order

Intuition: for most of the part, they are the same, but for one (some "in  $F_{\text{no}}$ ) place that is "controllable different".

Smooth

---

But: if we would like to lower bound:

$$\mathbb{E} \left[ \int (\hat{m}(x) - m(x))^2 dx \right]$$

then a Le Cam method would not have a tight bound. We will use more advanced methods.

---

Recap:  $R_n = \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}[d(\hat{\theta}, \theta(P))]$   
 $\geq \frac{5}{2} \cdot \inf_{\psi} \max_j P_j(\psi \neq j)$

We want  $L_n \leq R_n \leq U_n$

## 2. Fano's Method

Theorem 11.

$$\text{(Fano's Inequality)} \quad \max_j P_j(\psi \neq j) \geq \frac{1}{N} \sum_{j=1}^N P_j(\psi \neq j)$$

$$\geq 1 - \frac{n\beta + \log 2}{\log N}, \quad \text{with}$$

$$\text{(No Proof)} \quad N > 2, \quad \beta = \max_{j \neq k} \text{KL}(P_j, P_k)$$

With Fano's Inequality, we have

$$\star R_n \geq \frac{S}{2} \cdot \left[ 1 - \frac{n\beta + \log 2}{\log N} \right] \geq \frac{S}{4}, \quad \text{if } \beta \leq \frac{\log N}{4n}$$

— "Fano Minimax Bound" (Corollary 13)

## 3. Tsybakov Bound

Theorem 14. For  $N \geq 3, P_j \ll P_0$ , if  $\frac{1}{N} \sum_{j=1}^N \text{KL}(P_j, P_0) \leq \frac{\log N}{16n}$ ,  
Given  $X_1, X_2, \dots, X_n \sim P \in \mathcal{P}$

then:

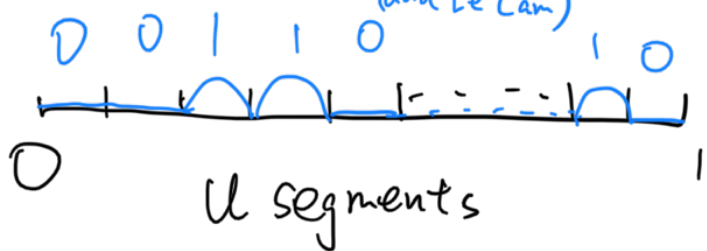
$$\star \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}[d(\hat{\theta}, \theta(P))] \geq \frac{S}{16}$$

where  $S = \min_{0 \leq j < k \leq N} d(\theta(P_j), \theta(P_k))$

△ Example 2: Lower bound  $\int_0^1 (\hat{m}(x) - m(x))^2 dx$ ,

where  $Y_i = m(X_i) + \epsilon_i$   $m \in \{m_0, m_1, \dots, m_N\}$

"One bump is not sufficiently valid" (and Le Cam)



1 → add a bump  
0 → add no bump

On  $m_i \leftrightarrow$  a  $\{0, 1\}^u$  sequences  $\mathbb{1}_i$

$\Rightarrow N = 2^u$

Define  $d(m_j, m_k) = \int (m_j(x) - m_k(x))^2 dx$   
 $= a(u)^2 \|\mathbb{1}_j - \mathbb{1}_k\|_1$

Of course, Le Cam does not work this time.

Note: if we choose all  $N = 2^u$  functions,  $\Rightarrow$  "Hamming distance"  $\geq$  small

then  $S = \min_{j \neq k} d(m_j, m_k) = a(u)^2$

$\Rightarrow R_n \geq \frac{S}{16n}$  is not sufficient to make sense.

Therefore, we have to choose a subset of all  $\{0,1\}^u$ , where each pair is far enough from each other and total number is large enough as well (or otherwise they are "easy to be distinguish" and again lead to a trivial lower risk bound).

△ Lemma 15: Varshamov - Gilbert's Lemma:

For  $N \geq 8$ , There exists a set

$$\{w^0, w^1, \dots, w^N\} \subset \{0,1\}^u \text{ such that}$$

$(w^0 = [0,0,\dots,0]^T \in \mathbb{R}^u)$

(i) "Constant Proportion Difference Property"

i.e.  $d(w^j, w^k) \geq \frac{u}{8}$

(ii) "Exponential Number Hypothesis Property"

i.e.  $N \geq 2^{\frac{u}{8}}$ .

Proof: Random Sample and Hoeffding Inequality

Hoeffding:  $P(|S| \geq t) \leq 2e^{-2t^2/n}$



$$\Pr[|S_n - \mathbb{E}[S_n]| \geq t] \leq 2 \exp\left(-\frac{t^2}{\sum (b_i - a_i)^2}\right), t > 0$$

$$\text{or } \Pr[S_n - \mathbb{E}[S_n] \geq t] \leq \exp\left(-\frac{2t^2}{\sum (b_i - a_i)^2}\right), t > 0$$

$$\Rightarrow \Pr\left[d(w^i, w^j) < \frac{u}{8}\right] \leq \exp\left(-\frac{2 \cdot \left(\frac{3u}{8}\right)^2}{u}\right) \\ = \exp\left(-\frac{9}{32}u\right), \forall i, j \in [N]$$

$$\Rightarrow \Pr\left[d(w^i, w^j) \geq \frac{u}{8}, \forall i, j \in [N]\right] \geq 1 - C_N^2 \cdot \exp\left(-\frac{9}{32}u\right) \\ \geq 1 - \frac{(2^u)^2}{2} \cdot \frac{1}{e^{\frac{9}{32}u}} \\ \geq 1 - \frac{1}{e^{\frac{9}{32}u}}$$

Total number of combinations:

$$(2^u)^N$$

$$\Rightarrow \# \text{ satisfying combination} \geq (2^u)^N \cdot \left(1 - \frac{1}{e^{\frac{9}{32}u}}\right) > 1$$

$\Rightarrow \exists$  a combination satisfying

$$\begin{cases} d(w_i, w_j) \geq \frac{u}{8} \\ N \geq 2^{\frac{u}{8}} \end{cases}$$

□

With this V-G theorem, we may construct a set  $M$  with  $N = 2^{\frac{u}{8}}$  and  $S = \frac{u}{8}$ .

Example 2 continue:

$$Y_i = m(X_i) + \epsilon_i$$

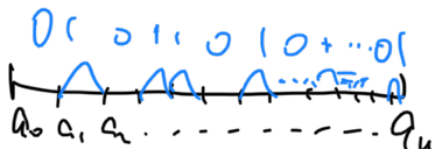
(. . .) (n-1) (β-1) . . .

$$\mathcal{M} = \{m: |m^{(k)}(x) - m^{(k)}(y)| \leq L \cdot |x-y|^k\}$$

( $\beta$ -order-smoothness family)

$$d(m_i, m_j) = \sqrt{\int (m_i - m_j)^2}$$

\* new definition  
Cor  $\|m_i - m_j\|_2$



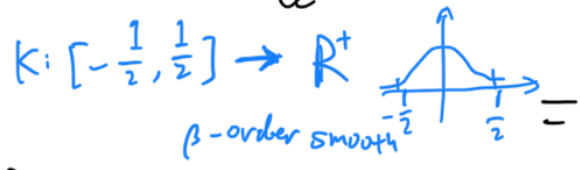
$u$  segments in total

$$= \sqrt{\left(\int (\text{one bump})^2\right)} \cdot \sqrt{\|1_i - 1_j\|_1}$$

each bump  $B_j(x)$

$$= \frac{L}{u^\beta} \cdot K\left(\frac{x - a_j}{\frac{1}{u}}\right)$$

$$\geq \frac{C}{u^{\beta + \frac{1}{2}}} \cdot \sqrt{\frac{u}{8}}$$



$$= C_1 \cdot u^{-\beta}$$

Also,  $KL(P_i, P_j) = C_2 \cdot u^{-2\beta}$

We omit the calculation

If we want to use Fano's method

we have to make  $KL(P_i, P_j) \leq \frac{\log N}{4n}$

$$= \frac{\log 2^{\frac{u}{8}}}{4n}$$

$$= \frac{u}{32n}$$

$$\Rightarrow C_2 \cdot u^{-2\beta} \leq \frac{u}{32n}$$

$$\Rightarrow U = (32 C_2 \cdot n)^{\frac{1}{2\beta+1}}$$

$$\begin{aligned} \Rightarrow R_n &\geq \frac{S}{4} = \frac{1}{4} \cdot C_1 \cdot U^{-\beta} \\ &= \Omega \left( n^{-\frac{\beta}{2\beta+1}} \right) \end{aligned}$$

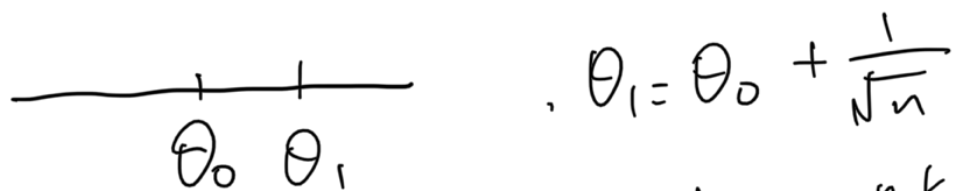
Note: for  $d$ -dimension setting, we have

$R_n \geq \Omega \left( n^{-\frac{\beta}{2\beta+d}} \right)$  in similar way.  
This matches the realized upper bound.

Example 4: (not to discuss)

Parametric model:  $R_n = \Theta \left( \frac{1}{\sqrt{n}} \right)$

i.e.,  $\inf_{\theta} \sup_{\hat{\theta}} \mathbb{E} |\hat{\theta} - \theta| \geq \frac{1}{\sqrt{n}}$



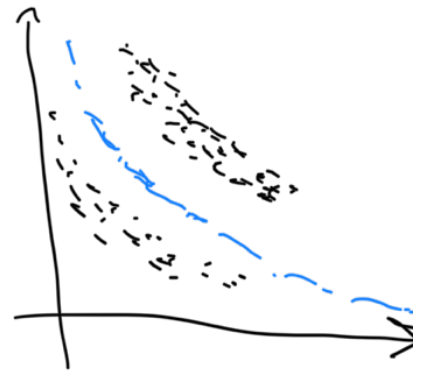
$$\theta_0 \quad \theta_1 \quad , \quad \theta_1 = \theta_0 + \frac{1}{\sqrt{n}}$$

then  $R_n \geq \frac{1}{\sqrt{n}} e^{-n \text{KL}(P_{\theta_0}, P_{\theta_1})}$  by LeCam

Recap:  $\text{KL}(P_0, P_1) \approx \frac{1}{n} \cdot \mathbb{I}(\theta)$

# Fisher Information

Example 5: (not to taught)  
semi-supervised learning  
e.g. no-label classification



Example 6:  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} P$  (if time permitted)

Usually:  $R_n = \Theta\left(n^{-\frac{2\beta}{2\beta+1}}\right)$  (for square loss)

but:

Observe:  $Z_i = X_i + \varepsilon_i$ ,  $\varepsilon_i \sim N(0, \sigma^2)$  but no access to  $X_i$

Estimate the density of  $X_i$

$\Rightarrow R_n = \Theta\left(\left(\frac{1}{\log n}\right)^\beta\right)$ , much slower!

Similarly, if we'd like to estimate  $m(\cdot)$  from  $Y_i = m(X_i) + \varepsilon_i$  but only observe

1  $Z_i = X_i + \varepsilon_i \sim N(m(X_i), \sigma^2)$ , then the  $R_n$ 's also logarithmic. i.e. we basically cannot estimate  $m(\cdot)$

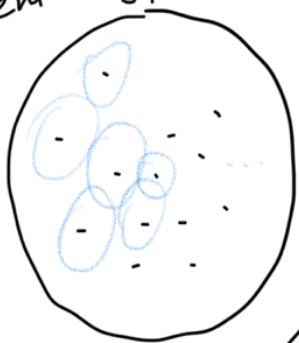
Very surprising!

(Note: if we only want  $Z_i \rightarrow Y_i$  prediction, then it's still fine, but it's not sth that fits  $m(\cdot)$  well.)

Difference between prediction and function estimating.

In general,

problem hypothesis family



Distinguishability

concentration

likelihood

$\downarrow \arg \max$

Upper bound

discretization

covering

$\downarrow$  count covering number

Lower bound  $N(\epsilon)$

Le Cam equation

$$\log N(\epsilon) = n \epsilon^2 \Rightarrow \epsilon_n \text{ is exactly the minimax}$$

↓  
"Entropy"

under some certain  
condition.