

Minimax method for evaluation

"Have you try your best?" This kind of question is often philosophical and mathematical, which would constantly bother you if you wonder the possibility to improve again.

In statistics, for a given problem, we usually need to impose a suitable model framework and solve it by any well-defined statistics of observed data, which leads to the same problem:

"Are you giving the best statistics?" To evaluate this, it's crucial to know what is "the best".

Example 1

Problem: the mean of population μ

Model: $\mu \in \mathbb{R}$

Observed Data: $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$

Statistics: $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$

"Is \bar{X}_n the best?"

By Lehmann-Scheffe Theorem, we know for all unbiased $\hat{\mu}$, i.e. $E\hat{\mu} = \mu$, we have

$$Var_{\mu} \bar{X}_n \leq Var_{\mu} \hat{\mu},$$

which means that \bar{X}_n is optimal in some sense.

What if we know $\mu \leq 0$, indeed, $\min\{0, \bar{X}_n\}$ is always better than \bar{X}_n since it utilizes the prior information of μ , regardless of the fact that

$$\begin{aligned} E \min\{0, \bar{X}_n\} &= E \bar{X}_n \mathbf{1}(\bar{X}_n \leq 0) \\ &< E \bar{X}_n, \end{aligned}$$

which means that the unbiasedness may not be appropriate to be used in this situation.

Note that

$$\begin{aligned} &E |\min\{0, \bar{X}_n\} - \mu| \quad (\mu \leq 0) \\ &= E |\min\{-\mu, \bar{X}_n - \mu\}| \leq E |\bar{X}_n - \mu|, \end{aligned}$$

which implies that $E | \cdot - \mu |$ may be intuitive correct for evaluation of $\hat{\mu}$.

Decision Theory

Decision space: Θ , $\theta_0 \in \Theta$

Observed data: $X_1, \dots, X_n \sim (\mathcal{X}', P_{\theta})$

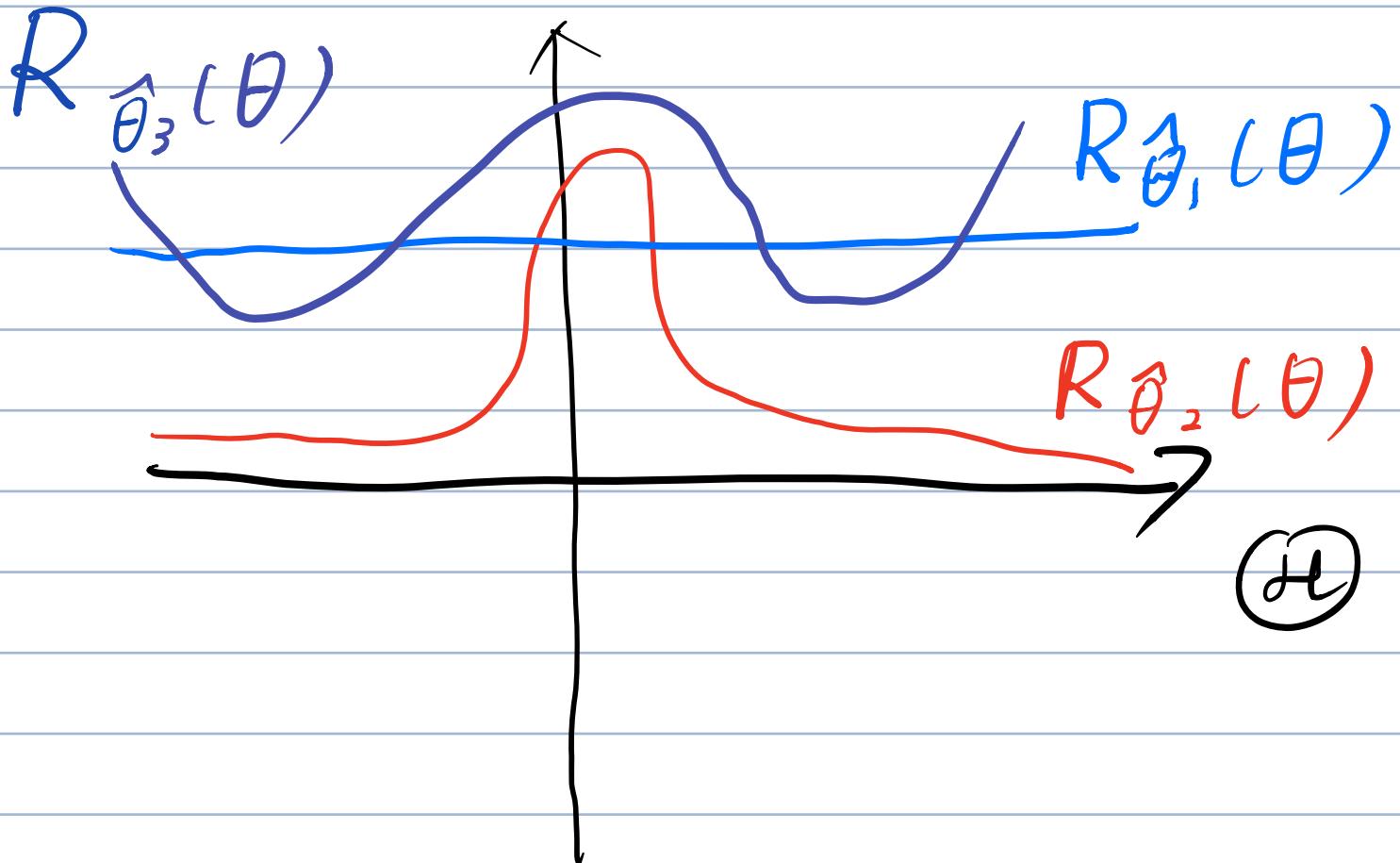
Target: θ_0 (density estimation)

Decision: $\hat{\theta}: \mathcal{X}' \rightarrow \Theta$

Loss function: $L: \Theta \times \Theta \rightarrow \mathbb{R}^+$

Risk function:

$$R_{\hat{\theta}}(\theta_0) = \mathbb{E}_{\theta_0} L(\hat{\theta}, \theta_0)$$



"Is $\hat{\theta}_1$ better than $\hat{\theta}_2$?!"

"Is $\hat{\theta}_3$ admissible?!"

① Average case: If we have some measures of Θ : π ,

$$\int R_\theta(\theta) \pi(d\theta)$$

could be used in comparison, where π is called the prior of Θ .

For the blessing of prior, the optimal $\hat{\theta}$:

$$\operatorname{argmin}_\theta \int R_\theta(\theta) \pi(d\theta)$$

exists under some mild conditions, which is the functionals of posterior distribution: $\pi(\theta | X_{1:n})$. For this reason, optimal decision is often called Bayes procedure.

② The worst case:

$$\sup_{\theta \in \Theta} R_\theta(\theta) \quad ①$$

we call the $\hat{\theta}$ minimizes ① the minimax optimal decision. (admissible?)

Minimal Optimality

$$M(\Theta) := \inf_{\hat{\theta}} \sup_{\theta \in \Theta} R_{\hat{\theta}}(\theta)$$

How to know $\hat{\theta}_0$ is minimax optimal?

① Upper bound of $R_{\hat{\theta}_0}(\theta)$:

$$U_{\hat{\theta}_0}(\Theta)$$

② Lower bound of $M(\Theta)$:

$$L(\Theta)$$

If $U_{\hat{\theta}_0}(\Theta) = L(\Theta)$, then

$$\begin{aligned} \sup_{\theta \in \Theta} R_{\hat{\theta}_0}(\theta) &\leq U_{\hat{\theta}_0}(\Theta) \\ &= L(\Theta) \leq M(\Theta) \end{aligned}$$

$\Rightarrow \hat{\theta}_0$ is minimax optimal.

Somehow, " $U_{\hat{\theta}_0}(\Theta) = L(\Theta)$ " is very strong, we can ease it:

$$U_{\hat{\theta}_0}(\Theta) \asymp L(\Theta) \quad ②$$

Remark: To achieve ②, $U_{\hat{\theta}_0}(\Theta)$ should be small enough, while $L(\Theta)$ should be large enough.

Example 2

In example 1, it implies that $M(\Theta) = 6^2/n$, if we constraint the decision $\hat{\mu}$ for μ should be unbiased, where

$$\frac{6^2}{n} = \sup_{\mu \in \Theta} \mathbb{E} (\bar{X}_n - \Theta)^2.$$

General Reduction Scheme

A Let $L = \Phi \circ d$, where $\Phi: \mathbb{R}^+ \rightarrow \mathbb{R}^+$ is an increasing function and d is a semi-metric of Θ , for a decision $\hat{\theta}_0$ of θ , if

$$\sup_{\theta \in \Theta} R_{\hat{\theta}_0}(\theta) = \sup_{\theta \in \Theta} \mathbb{E} \Phi(d(\hat{\theta}_0, \theta))$$

$$\leq \Phi(8),$$

simultaneously, note that

$$R_{\hat{\theta}}(\theta) = \mathbb{E}_{\theta} \mathbb{I}(d(\hat{\theta}, \theta))$$

$$\geq \underline{\Phi}(\delta) P_{\theta}(d(\hat{\theta}, \theta) > \delta)$$

if

$$\inf_{\theta} \sup_{\theta' \in \Theta} P_{\theta}(d(\theta, \hat{\theta}) > \delta) > 0$$

$$\Rightarrow M(\hat{\theta}) \geq \underline{\Phi}(\delta)$$

$\Rightarrow M(\hat{\theta}) \asymp \underline{\Phi}(\delta)$ and $\hat{\theta}_0$ is minimax optimal.

So what we need to do is to give a sharpen upper bound: $\underline{\Phi}(\delta)$, and then prove that

(3)

$$\inf_{\theta} \sup_{\theta' \in \Theta} P_{\theta}(d(\theta, \hat{\theta}) > \delta) > 0$$

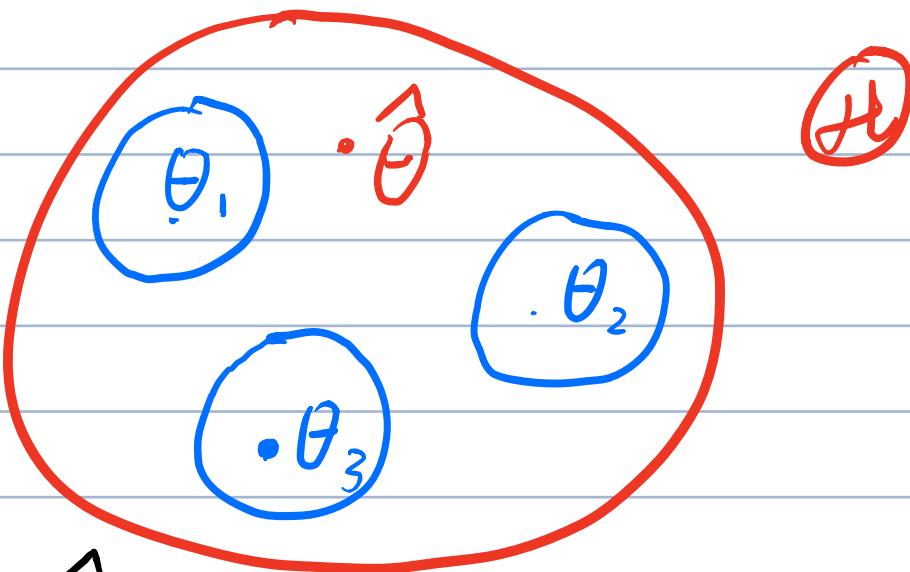
Remark:

$$\begin{aligned} \underline{\Phi} &\leq \frac{M(\hat{\theta})}{\underline{\Phi}(\delta)} \leq \frac{\sup_{\theta' \in \Theta} \mathbb{E} \mathbb{I}(d(\hat{\theta}, \theta))}{\underline{\Phi}(\delta)} \\ &\leq C. \end{aligned}$$

12 From estimation to testing

For given $\delta > 0$, to prove ②, take 2 δ -packing centers $\{\hat{\theta}\}_{i=1:M}$ of \mathcal{H}
 then: (If they exists!)

$$\{d(\hat{\theta}, \theta_i) < \delta\} \subset \bigcap_{j \neq i} \{d(\hat{\theta}, \theta_j) \geq \delta\}$$



define $\hat{\phi} \in \{1, \dots, M\}$:

$$\hat{\phi} := \phi(\hat{\theta}) = \operatorname{argmin}_j d(\hat{\theta}, \theta_j)$$

$$\Rightarrow \{d(\hat{\theta}, \theta_i) < \delta\} \subset \{\hat{\phi} = i\}$$

$$\Rightarrow P_{\theta_i}(d(\hat{\theta}, \theta_i) \geq \delta) \geq P_{\theta_i}(\hat{\phi} \neq i)$$

\Rightarrow

$$\begin{aligned}
 & \sup_{\theta \in \Theta} P_\theta(d(\hat{\theta}, \theta) > \delta) \\
 & \geq \sup_i P_{\theta_i}(d(\hat{\theta}, \theta_i) > \delta) \\
 & \geq \sup_i P_{\theta_i}(\hat{\phi} \neq i) \\
 & \geq \frac{1}{m} \sum_{i=1}^m P_{\theta_i}(\hat{\phi} \neq i)
 \end{aligned}$$

$$\begin{aligned}
 \Rightarrow & \inf_{\hat{\phi}} \sup_{\theta \in \Theta} P_\theta(d(\hat{\theta}, \theta) > \delta) \\
 & \geq \inf_{\hat{\phi}} \frac{1}{m} \sum_{i=1}^m P_{\theta_i}(\hat{\phi} \neq i),
 \end{aligned}$$

where $\hat{\phi}$ is any statistics valued in $\{1, \dots, M\}$. (fixed given data)

Then given $\delta > 0$, we just know to prove that $\exists \{\theta_i\}_{i=1:m}, d(\theta_i, \theta_j) \geq 2\delta$

s.t.

$$\inf_{\hat{\phi}} \frac{1}{m} \sum_{i=1}^m P_{\theta_i}(\hat{\phi} \neq i) > 0.$$

Le Cam's method

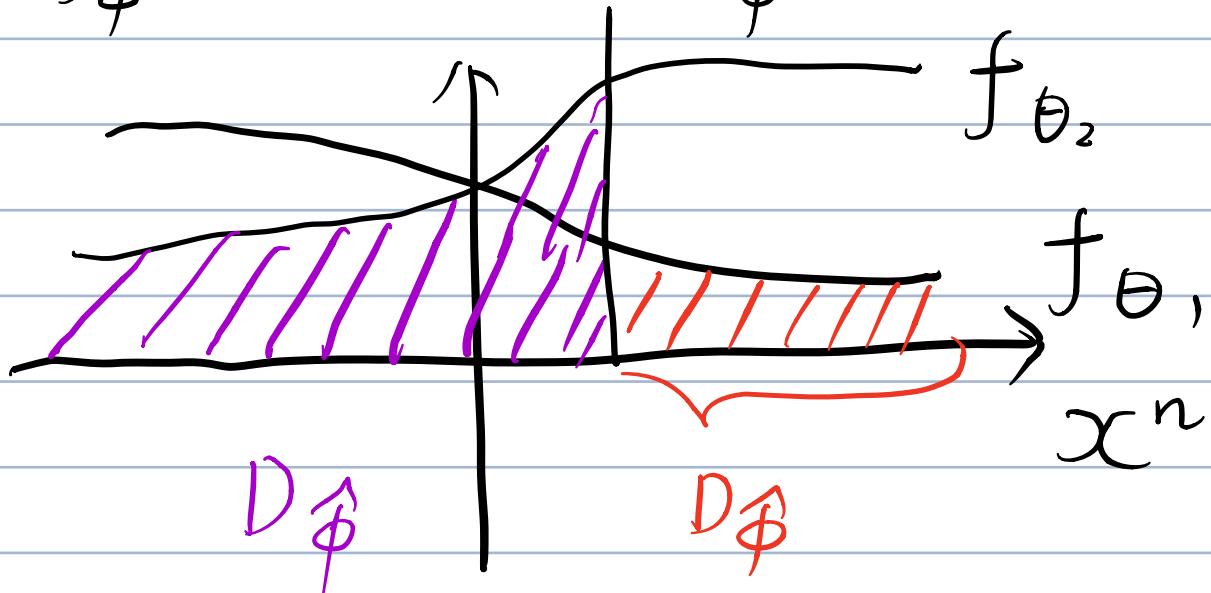
If $M=2$, we need to prove

$$\inf_{\hat{\phi}} \frac{1}{2} (P_{\theta_1}(C_{\hat{\phi}} \neq 1) + P_{\theta_2}(C_{\hat{\phi}} \neq 2)) \geq 0$$

We assume that $P_{\theta_1} < v$, and $D_{\hat{\phi}} = \{x^n \in X^n : \hat{\phi}(x^n) \neq 1\}$, then

$$P_{\theta_1}(C_{\hat{\phi}}) + P_{\theta_2}(C_{\hat{\phi}})$$

$$= \int_{D_{\hat{\phi}}} f_{\theta_1} dv + \int_{D_{\hat{\phi}}^c} f_{\theta_2} dv$$



$$\geq \int \min\{f_{\theta_1}, f_{\theta_2}\} dv,$$

"=" achieve iff $D_{\hat{\phi}} = \{x^n : f_{\theta_1} \leq f_{\theta_2}\}$.

$$\Rightarrow \inf_{\hat{\phi}} P_{\theta_1}(D_{\hat{\phi}}) + P_{\theta_2}(D_{\hat{\phi}}^c)$$

$$\geq \int \min \{f_{\theta_1}, f_{\theta_2}\}$$

$$= 1 - \frac{1}{2} \int |f_{\theta_1} - f_{\theta_2}|$$

$$= 1 - \|P_{\theta_1} - P_{\theta_2}\|_{TV}, \text{ where}$$

$$\|P_{\theta_1} - P_{\theta_2}\| = \sup_A |P_{\theta_1}(A) - P_{\theta_2}(A)|.$$

Lemma 1 (Neyman-Pearson)

$$\inf_{\hat{\phi}} P_{\theta_1}(D_{\hat{\phi}} \neq 1) + P_{\theta_2}(D_{\hat{\phi}} \neq 2)$$

$$\geq 1 - \|P_{\theta_1} - P_{\theta_2}\|_{TV},$$

where " $=$ " achieve iff likelihood ratio

$$\hat{\phi} = \begin{cases} 1, & \underline{f_{\theta_1} \geq f_{\theta_2}} \\ 2, & \text{else} \end{cases}$$

Remark: Indeed, $d(\theta_1, \theta_2) \geq 28$
 and $1 - \|P_{\theta_1} - P_{\theta_2}\|_{TV} > 0$ should be
 achieved both, which is a trade-off
 for suitable $\theta_1, \theta_2 \in \Theta$).

Hellinger distance

Let P, Q be two probability measures of X^n ($P, Q \ll \nu$),

$$\begin{aligned} \mathcal{H}(P, Q) &:= \left(\int (\sqrt{P} - \sqrt{Q})^2 \right)^{\frac{1}{2}} \\ &= (2 - 2 \int \sqrt{PQ})^{\frac{1}{2}} \end{aligned}$$

Lemma 2

$$\frac{(\mathcal{H}(P, Q))^2}{2} \leq \|P - Q\|_{TV} \leq \mathcal{H}(P, Q)$$

Pf: $\frac{(\mathcal{H}(P, Q))^2}{2} = 1 - \int \sqrt{PQ}$
 $\leq 1 - \int \min\{P, Q\} = \|P - Q\|_{TV}$

$$\begin{aligned} \mathcal{H}(P, Q) &= 2 - 2 \int \sqrt{PQ} \\ &= 2 - \int \sqrt{\max\{P, Q\} \min\{P, Q\}} \\ &\geq 2 - \left(\int \max\{P, Q\} \right)^{\frac{1}{2}} \left(\int \min\{P, Q\} \right)^{\frac{1}{2}} \\ &= 2 - 2 \left(1 + \|P - Q\|_{TV} \right)^{\frac{1}{2}} \left(1 - \|P - Q\|_{TV} \right)^{\frac{1}{2}} \end{aligned}$$

$$= 2 - 2 C (1 - \|P - Q\|_{TV}^2)^{\frac{1}{2}}$$

$\geq \|P - Q\|_{TV}$ since

$$2 - 2C(1-x)^{\frac{1}{2}} \geq x, 0 \leq x \leq 1$$

□

Kullback-Leibler divergence

$$KL(P||Q) = \int p \log \frac{p}{q}$$

Lemma 2

$$KL(P||Q) \geq (\mathcal{H}(P, Q))^2$$

Pf:

$$KL(P||Q) = \int p \log \frac{p}{q}$$

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

$$\geq -2 \int p \left(\sqrt{\frac{q}{p}} - 1 \right)$$

$$= 2 - 2 \int \sqrt{pq}.$$

□

Decoupling Property

Now we mark $P^n = P \otimes \dots \otimes P$
and P is a probability measure of X .

Pro

$$\mathcal{H}^2(P^n, Q^n) \leq n \mathcal{H}^2(P, Q)$$
$$KL(P^n || Q^n) = n KL(P || Q)$$

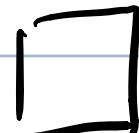
Pf:

$$\frac{1}{2} \mathcal{H}^2(P^n, Q^n) = 1 - \int \sqrt{P^n Q^n}$$

$$= 1 - (\int \sqrt{pq})^2 = 1 - (1 - \frac{1}{2} \mathcal{H}^2(P, Q))^n$$

$$\leq \frac{1}{2} n \mathcal{H}^2(P, Q) \quad (1 - (1-x)^n \leq nx)$$

$$KL(P^n || Q^n) = \int P^n \log \frac{P^n}{Q^n}$$
$$= \int_{i=1}^n \int p(x_1, \dots, x_n) \log \frac{p(x_i)}{q(x_i)} dx_i$$
$$= n \quad KL(P || Q)$$



Example 3

In example 1, we know

$$KLC(P_{\mu_1}, P_{\mu_2}) = \frac{(\mu_1 - \mu_2)^2}{26^2},$$

then

$$1 - \|P_{\mu_1} - P_{\mu_2}\|_{TV} \geq 1 - n \frac{(\mu_1 - \mu_2)^2}{26^2},$$

then we should choose δ s.t.

$$1 - \frac{(\mu_1 - \mu_2)^2}{26^2/n} > 0 \quad \text{and}$$

$$|\mu_1 - \mu_2| \geq 28,$$

$$\Rightarrow \delta \leq \sqrt{\frac{6^2}{n}}$$

$$\Rightarrow M(\mathcal{H}) \geq \Phi(\sqrt{\frac{6^2}{n}})$$

If $\Phi(t) = t^2$, then

$$\sup_{\mu} E(\bar{x}_n - \mu)^2 = \frac{6^2}{n}$$

$\Rightarrow \bar{x}_n$ is minimal optimal.

Example 4

$X_1, \dots, X_n \stackrel{iid}{\sim} U[\theta, \theta+1]$, we use $\min\{X_i\}$ to estimate θ :

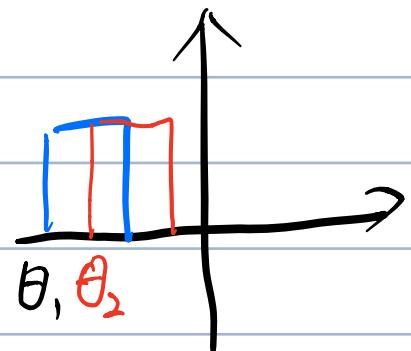
$$\begin{aligned} & \mathbb{E} |\min\{X_i\} - \theta| \\ &= \mathbb{E} \min\{Y_i\}, Y_i \stackrel{iid}{\sim} U[0, 1] \\ &= \int \mathbb{P}(\min\{Y_i\} > t) dt \\ &= \frac{1}{n+1} \end{aligned}$$

And $H^2(CU[\theta, \theta+1], CU[\theta_2, \theta_2+1])$

$$\begin{aligned} &= \begin{cases} 2|\theta - \theta_2|, & \text{if } |\theta - \theta_2| \leq 1 \\ 2, & \text{else} \end{cases} \\ &= 2 \min\{|\theta - \theta_2|, 1\} \end{aligned}$$

$$\Rightarrow H^2(P_\theta, P_{\theta_2}) \leq n \min\{|\theta - \theta_2|, 1\}$$

$$\Rightarrow \delta = \frac{1}{2n}$$



$$\Rightarrow M(\theta) \geq \bar{\Phi}(1/n)$$

$\Rightarrow \min\{X_i\}$ is minimax optimal.

Le Cam for functional

Compared with the traditional way to evaluate a point estimator, such as, **UMVUE** (Uniformly minimal variance unbiased estimator) or **Cramér-Rao Lower bound**, minimax method is less powerful and even boring for the parametric model. But the situation gets little different when we discuss the point estimator for **non-parametric** model, since the likelihood of data has an "ambiguous" form of parameter (**principle of sufficiency and likelihood is hard to be used**) and the procedure of **model selection** makes the evaluation getting hard, while the minimax aspect mostly focuses on the geometry structure of the parameter space, which is applicable for the evaluation of wide range of point estimator.

Now we let $\mathcal{F} := \mathbb{H}$, which are some density function spaces.

Similarly, we can construct some non-parametric statistics for $f \in \mathcal{F}$: f_n and evaluate it by

$$\sup_{f \in \mathcal{F}} \mathbb{E} \Phi_C d(f_n, f),$$

not that far, we firstly focus on some functional of $\mathbb{E} \Phi_C$ mean, median, or evaluation functional $\Theta(\cdot)$: θ , and given a statistic for $\Theta(f)$: $\hat{\theta}$, define

$$\inf_{\hat{\theta}} \hat{\theta} \sup_{f \in \mathcal{F}} \mathbb{E} \Phi_C d(\hat{\theta}, \Theta(f)) \\ := M(\Theta(\mathcal{F})).$$

The part gets harder

By Le Cam, if $\exists f, g \in \mathcal{F}$, s.t.

$$d(\Theta(f), \Theta(g)) \geq 2\delta,$$

for given $\delta > 0$, then

$$M(\Theta(\mathcal{F})) \geq \Phi(\delta)$$

if

$$1 - \|P_f^n - P_g^n\|_{TV} > 0.$$

Note that

$$\|P_f^n - P_g^n\|^2 \leq H^2(P_f^n, P_g^n)$$
$$\leq n H^2(P_f, P_g),$$

If we know that

$$d(\Theta(f), \Theta(g)) \geq CH(P_f, P_g) \quad \textcircled{4}$$

Then if $\exists f, g \in \mathcal{J}$, s.t.

$$H^2(P_f, P_g) = \frac{1}{4n}$$

\Rightarrow

$$M(\Theta(\mathcal{J})) \geq \Phi\left(\frac{c}{4n}\right)$$

If $\textcircled{4}$ is hard to achieve, we can define

$$W_\theta(\varepsilon) = \sup_{\substack{f, g \in \mathcal{J} \\ H(f, g) \leq \varepsilon}} |\Theta(f) - \Theta(g)|$$

$$\Rightarrow M(\Theta(\mathcal{J})) \geq \Phi\left(\frac{1}{2} W_\theta\left(\frac{1}{2n}\right)\right).$$

Example 5 (Minimax rate of pointwise estimation of Lipschitz density)

$$\bar{\mathcal{H}}_L = \{f \in C[0,1], |f(x) - f(y)| \leq L|x-y|\},$$

$$\int f(x) dx = 1 \text{ (with } \| \cdot \|_{\sup} \text{)}$$

There is a famous estimator for this kind of density: $\hat{f}_{n,\lambda}$, kernel density estimation, and

$$\sup_{f \in \bar{\mathcal{H}}_L} \mathbb{E}(\hat{f}_{n,\lambda}(x_0) - f(x_0))^2$$

$$\leq n^{-\frac{2}{3}}, x_0 \in [0,1],$$

if we choose a suitable window width $\lambda \propto \lambda_n$, the optimal λ depends on n).

Define $\Theta(f) = f(x_0)$, $f \in \bar{\mathcal{H}}_L$.

Now we let $L=1$, choose $f \in \bar{\mathcal{H}}_1$: $f \equiv 1$, and $\phi \in \mathcal{H}_1$, $\int \phi = 0$, then

$$g := f + \phi \in \bar{\mathcal{H}}_1,$$

$$\mathcal{H}^2(f, g) = 2 - 2 \int \sqrt{1 + \phi(x)} dx.$$

Let $\|f - g\|_{\sup} = \|\phi\|_{\sup} = 8$.

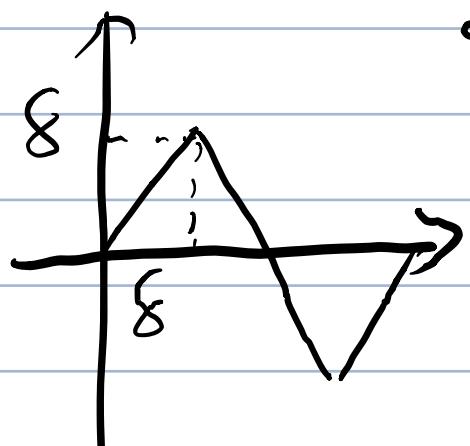
In addition, if $\phi(0) = 0$, then

$$\begin{aligned} \frac{1}{2} \mathcal{H}^2(f, g) &= \int \sqrt{1 + \phi(x)} \Big|_x^0 dx \\ &\leq \int C_1 \phi(x) + C_2 \phi^2(x) dx, \end{aligned}$$

(where $C_1 = u'(0)$, $C_2 = \sup_{x \in S} u''(x)$,)
 $u(S) = \sqrt{1 + 8}$.

$$\leq \int \phi^2(x) \leq \int_0^8 t^2 dt \asymp 8^3$$

Let $8^3 \asymp \frac{1}{n}$



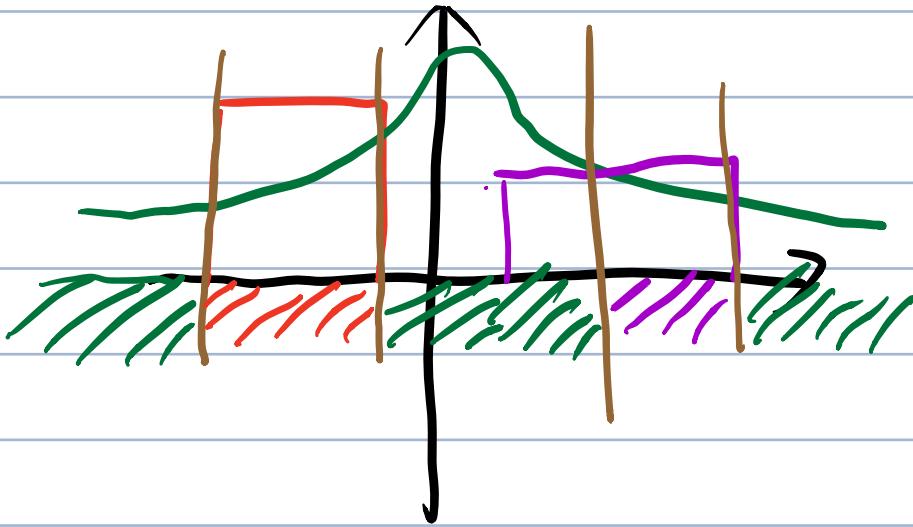
$$\Rightarrow W_0(\frac{1}{2n}) \gtrsim n^{-\frac{1}{3}}$$

$$\Rightarrow M(\Theta \mathcal{F}) \gtrsim \Theta(n^{-\frac{1}{3}}).$$

Fano's method

For $M \neq 2$, we need to evaluate the error of multiple testing, i.e.

$$\begin{aligned} & \frac{1}{M} \sum_{i=1}^M P_{\theta_i} (\hat{\phi} \neq i) \\ &= \frac{1}{M} \sum_{i=1}^M \int_{D_{\hat{\phi}}^i} f_{\theta_i} \\ &\geq \frac{1}{M} \int \min_i \{f_{\theta_i}\} \end{aligned}$$



An interesting problem is that "Would the optimal error get smaller and smaller if you make more and more comparisons?"

To prove this intuition, we need to introducing some elements of Information Theory.

Entropy

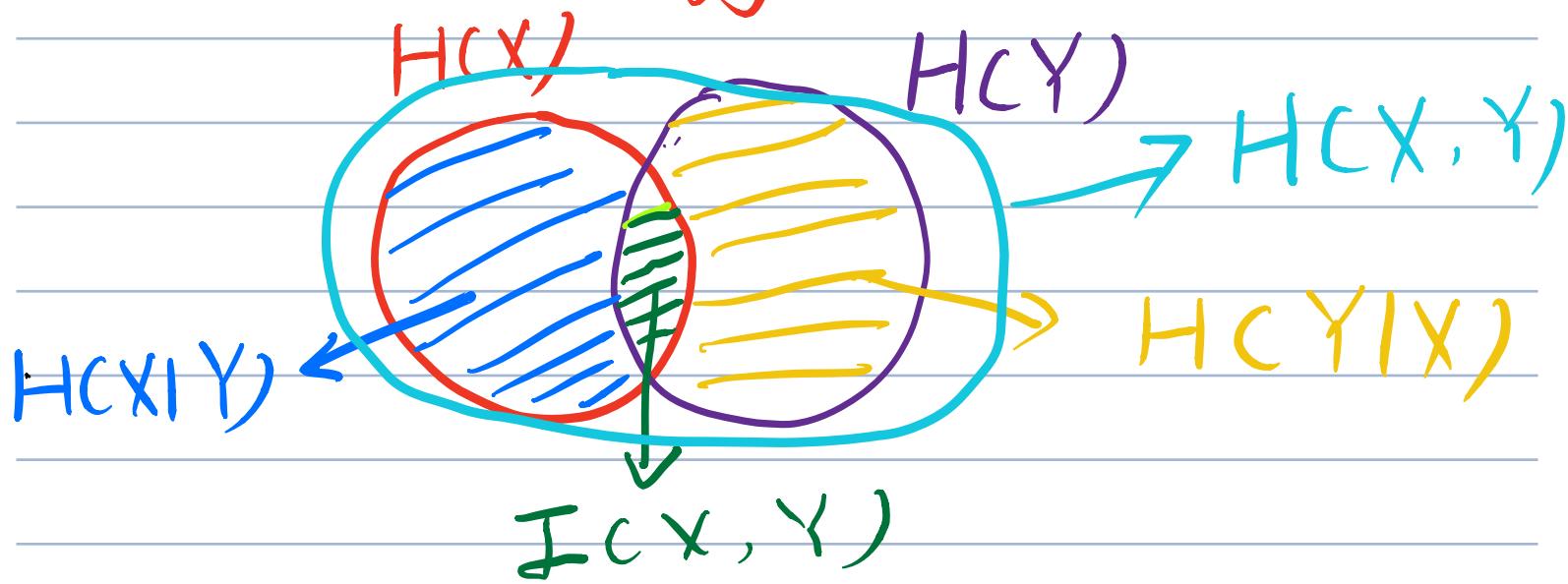
X, Y are two random elements (no need the same), which joint distribution exists.

$I(X; Y) := KL(P_{X,Y} \| P_X P_Y)$, which quantifies the dependence of X and Y .
(Mutual Information)

$H(X) := - \int p_x \log p_x$, which quantifies the information (randomness) of X .
(Shannon entropy)

$H(X|Y=y) := H(Z), Z \sim X|Y=y$.

$H(X|Y) = \int H(X|Y=y) P_Y(dy)$.
(Conditional entropy)



Remark:

① If X is a random element valued in $X, |X| < \infty$, then

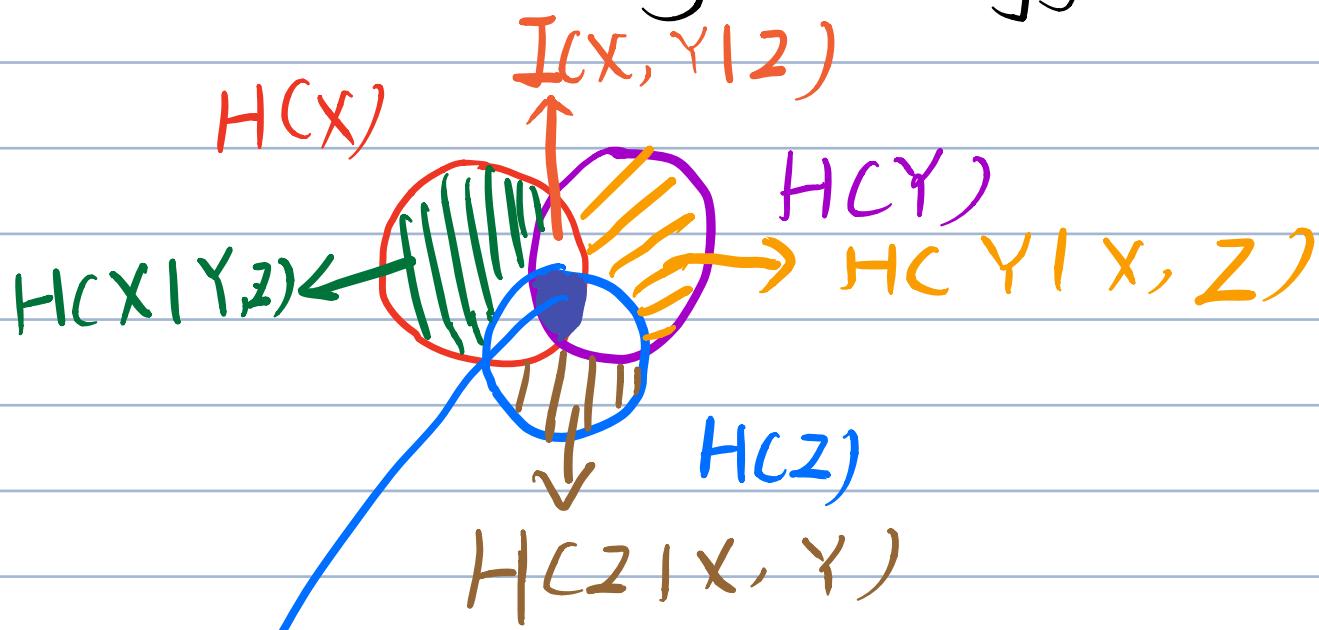
$$\begin{aligned} H(X, U) &= H(U) + H(X|U), \\ (\text{where } U = F^{-1}(X) \sim \text{Unif}\{1, \dots, |X|\}) \\ &= H(U) = \log |X| \end{aligned}$$

$$\Rightarrow H(X) \leq H(X, U) = \log |X|.$$

\Rightarrow the uniform distribution is the most random element.

② We can define $H(X, Y, Z)$ similarly:

$$H(X, Y, Z) = -\int p_{x,y,z} \log p_{x,y,z}.$$



Not well-defined! And $I(X, Y|Z)$ may be larger than $I(X, Y)$.

Return to our problem, let $U \cup U_{\text{inf}}$

g_1, \dots, g_M , $Z|U = i$, P_{θ_i} , then

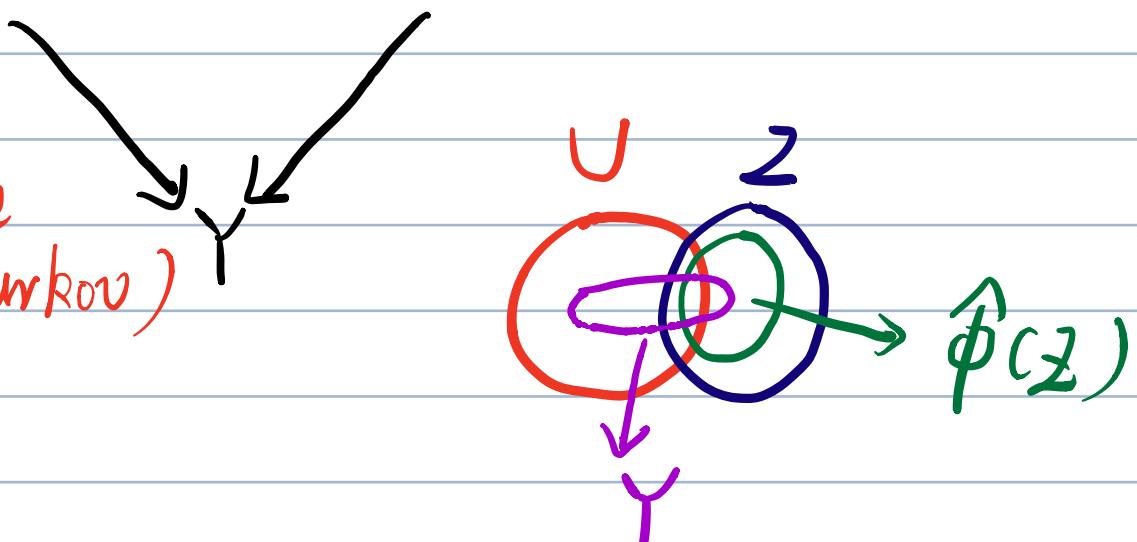
$$\frac{1}{M} \sum_{i=1}^M P_{\theta_i}(\hat{\phi} \neq i) = P(\hat{\phi}(z) \neq U) := p,$$

Let $Y = \mathbb{1}(\hat{\phi}(z) \neq U)$, then

$Y \sim BC(1, p)$.

$$U \rightarrow Z \rightarrow \hat{\phi}(z)$$

Information
transfers
through the
chain (Markov)



$$\Rightarrow H(U, Y | Z) = H(U | Z)$$

$$= H(U | Y, Z) + H(Y | Z)$$

$$= p H(U | Y=1, Z) + H(Y | Z)$$

$$\leq p \log(M-1) + H(Y)$$

$$= p \log(M-1) - p \log p - (1-p) \log(1-p)$$

$$\leq p \log(M-1) + \log 2$$

$$\Rightarrow p \geq \frac{H(U | Z) - \log 2}{\log(M-1)}$$

$$\Rightarrow p \geq \frac{\log M - I(U; Z) - \log^2}{\log(M-1)}$$

$$\geq 1 - \frac{I(U; Z) + \log^2}{\log M}.$$

(Fano)

If $\exists \{\theta_i\}_{i=1:M} \subset \Theta$ for given $\delta > 0$ s.t.
 $d(\theta_i, \theta_j) \geq 2\delta$, and

$$1 - \frac{I(U; Z) + \log^2}{\log M} > 0, \text{ then}$$

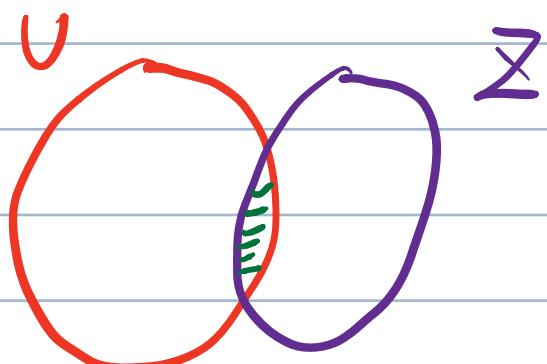
$I(U; Z) \geq \underline{I}(\delta)$, where
 $U \sim \text{Unif}\{1, \dots, M\}$, $Z|U=i \sim P_{\theta_i}$.

Remark:

① $\log M = H(U) = I(U; U)$, as $\delta \rightarrow 0^+$,
 if

$$\frac{I(U; Z)}{I(U; U)} \rightarrow 0^+,$$

then Fano holds.



$$\begin{aligned}
 ② I(U; Z) &= KLC P_{U,Z} || P_U P_Z \\
 &= KLC P_{Z|U} P_U || P_Z P_U \\
 &= \int P_{Z|U} P_U \log \frac{P_{Z|U}}{P_Z} \\
 &= \frac{1}{M} \sum_{i=1}^M I_{\theta_i} \log \frac{P_{\theta_i}}{P_Z} \\
 &= \frac{1}{M} \sum_{i=1}^M KLC P_{\theta_i} || \frac{1}{M} \sum_{k=1}^M P_{\theta_k} \\
 &\leq \frac{1}{M^2} \sum_{i,k} KLC (P_{\theta_i} || P_{\theta_k})
 \end{aligned}$$

Example 6 (Minimax risk for linear regression)

$Y = X\beta_0 + \epsilon$ X is fixed, and $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2 I_n)$, we want to evaluate

$$M(\Theta) := \inf_{\beta} \sup_{\mu \in \Theta} \mathbb{E} \|\hat{\mu} - \mu\|_n^2, \mu = X\beta$$

where $\Theta = \text{Im}(X)$. (Identifiable)

Let $r = \dim(\Theta) = \text{rank}(X)$, then

$$P(\mathcal{E}; B_\Theta(c_0; d)) \asymp \left(\frac{d}{\delta}\right)^r,$$

and

$$\begin{aligned}
 KLC P_{\beta} || P_{\beta_0} &= \frac{\|X(\beta - \beta_0)\|_2^2}{2\sigma^2} \\
 &= \frac{\|\mu_1 - \mu_2\|_n^2}{2\sigma^2} \leq \frac{n\delta^2}{6^2}
 \end{aligned}$$

To ensure $0 < \frac{\delta^2 / 6^2}{r \log d / \delta} < 1$

Let $d = 28$, $\delta = \sqrt{\frac{2}{n}} 6$, then
 $M(\Theta) = \Phi c \sqrt{\frac{2}{n}} 6$.

Example 7 (Minimax risk for Sparse mean)

(Gibert - Varshamov)

$\mathcal{A} := \{ \theta \in \{0,1\}^P, \|\theta\|_0 = s \}$, if $1 \leq s \leq P/8$, then there exists $\{\theta_i\}_{i=1}^M \in \mathcal{A}$ s.t.

$$\log M \geq s \underbrace{\log(1 + \frac{P}{2s})}$$

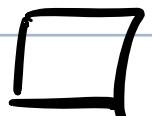
and

$$\|\theta_i - \theta_j\|_0 \geq s/2, i \neq j.$$

Pf: let $X_i \stackrel{iid}{\sim} \text{Unif}(\mathcal{A})$, $i = 1, \dots, M$
 we want to prove

$$P\left(\bigcap_{i \neq j} \|X_i - X_j\|_0 \geq s/2\right) > 0 \text{ if}$$

$$\log M \geq s \log(1 + \frac{P}{2s})$$



Now we have

$Y = \mu + \varepsilon$, $\beta_0 \in \mathbb{R}^P$ and
 $\|\mu\|_0 \leq s$, $\varepsilon \sim N(0, 6^2 I_n)$, similarly:

$$KL(P_\mu^n, P_{\mu_2}^n) = n \|\mu_1 - \mu_2\|^2 / 26^2$$

In order to ensure, let $\|\mu_i - \mu_j\| \geq 28$,
 in order to ensure

$$0 < \frac{n(28)^2 / 26^2}{\log M} < 1,$$

Let $\delta \asymp \sqrt{s \log(1 + \frac{P}{2s})}/n \cdot 6$, indeed
if $\mu_i \neq \mu_j$ s.t. $\|\mu_i - \mu_j\| \geq 2\delta$ exists since

Let $\mu_i = \theta_i \delta$ is, then

$$\|\mu_i - \mu_j\| = \frac{\delta}{s} \|\theta_i - \theta_j\|_0 \geq \delta.$$