

Chapter 3. M-estimator based on Le Cam

Setup: $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} P_{\theta_0}$ ($=: P_0$)
the true DG Parameter.
 $\Theta = \{\theta \in \mathbb{R}^k : \theta_0 \in \Theta\}$

$\Phi: \Theta \mapsto \mathbb{R}^K$ e.g. population mean
Goal of interest: $\phi_0 = \Phi(\theta_0)$

M-estimator Framework

① $\text{Im}(\Phi) \subset S$

↑
Image

② Restriction on Φ :

$$\forall \theta \in \Theta, \Phi(\theta) \in \underset{\phi \in S}{\arg \max} E_\theta [m_\phi(x)]$$

$x \sim P_\theta$

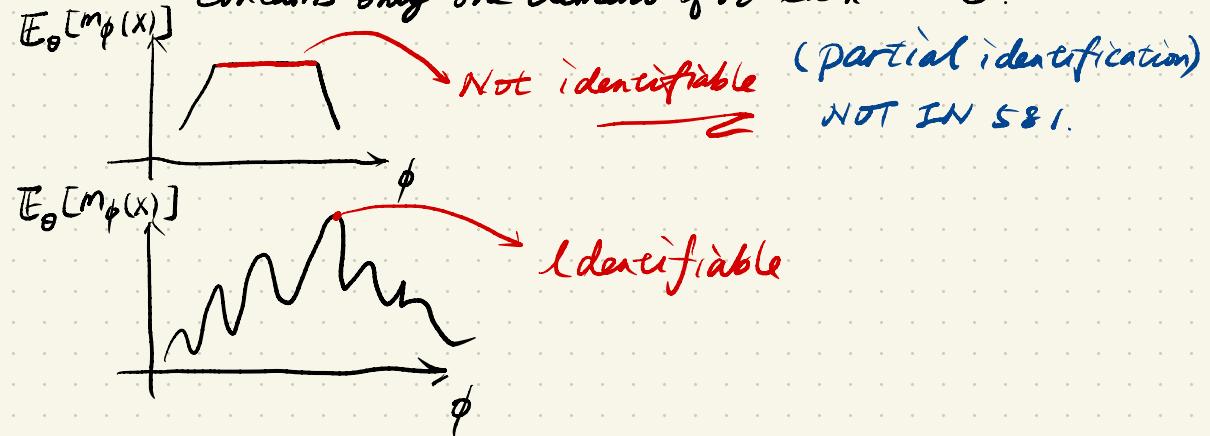
where m_ϕ can be understood as

- (a) negative of a loss function;
- (b) **pseudo** likelihood function;

M-est inference

① Identification

" m_ϕ " identifies $\Phi(\theta)$ iff the argmax contains only one element for each $\theta \in \Theta$.



② Regarding estimation:

Recall $\hat{\Phi}(\theta_0) = \underset{\phi \in S}{\operatorname{argmax}} E_{\theta_0}[m_\phi(x)]$

$$\hat{\theta}_n \in \underset{\phi}{\operatorname{argmax}} \hat{E}[m_\phi(x)]$$

\uparrow

m -estimator of

$$\hat{\Phi}(\theta_0)$$

$$= \underset{\phi}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n m_\phi(x_i)$$

Notation simplification

$$Pf := \int f(x) P(dx) (\mathbb{E}_{x \sim P})$$

P_n : the empirical probability measure

$$\text{e.g. } P_m \phi = \int m_\phi(x) P(dx)$$

$$P_0 m_\phi = \int m_\phi(x) P_{\theta_0}(dx)$$

$$P_n m_\phi = \frac{1}{n} \sum_{i=1}^n m_\phi(x_i)$$

$$\phi_0 \in \arg \max_{\phi} P_0 m_\phi$$

$$\phi_n \in \arg \max_{\phi} P_n m_\phi$$

$$\text{CLT: } \sqrt{n} (P_n - P_0) f \Rightarrow N(0, \Sigma_f)$$

Example (MLE as an M-est.)

KL-divergence: Given any P, Q , assuming

$$\begin{aligned} P \ll \mu \text{ and } Q \ll \mu &\rightarrow \text{Lebesgue/Counting measure usually.} \\ (\Rightarrow P = \frac{dP}{d\mu} \text{ and } Q = \frac{dQ}{d\mu}) \end{aligned}$$

$$D_{KL}(P || Q) := -P \left[\log \underbrace{\frac{Q}{P}}_{m_\theta} \right]$$

$$\theta_0 \in \arg \max_{\theta \in \Theta} [-D_{KL}(P_{\theta_0} || P_0)]$$

$$\theta_n \in \operatorname{argmax} P_n m_\theta \quad \xrightarrow{\text{log } \frac{q}{P}}$$

We can show, by Jensen \leq (Gibbs ineq) that

$$(1) D_{KL}(P||Q) \geq 0$$

$$(2) D_{KL}(P||Q) = 0 \text{ iff } P = Q.$$

$$\Downarrow$$

$$① \theta_0 \in \operatorname{argmax}_{\theta \in \Theta} [-D_{KL}(P_{\theta_0} || P_\theta)]$$

$$= \operatorname{argmax}_{\theta \in \Theta} P_{\theta_0} \left[\log \frac{P_\theta}{P_{\theta_0}} \right]$$

$$= \operatorname{argmax}_{\theta \in \Theta} P_{\theta_0} \left[\log P_\theta \right]$$

If $P_0 \log P_0 < \infty$

m_ϕ with $\Phi: \theta \mapsto \theta$.

$$= \operatorname{argmax}_{\theta \in \Theta} P_{\theta_0} m_\phi$$

② Identifiability: $m_\phi(m_\theta)$ identifies the M-est if

$$P_{\theta_1} = P_{\theta_2} \text{ iff } \theta_1 = \theta_2.$$

$$③ \theta_n \in \operatorname{argmax} P_n m_\theta$$

$$= \operatorname{argmax}_{\theta \in \Theta} P_n [\log P_\theta]$$

$$= \operatorname{argmax}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log P_\theta(X_i)$$

$$= \operatorname{argmax}_{\theta \in \Theta} \prod_{i=1}^n P_\theta(X_i)$$

Classic def. of MLE

Z -est (a cousin of M -est)

Setup: the same as M -est

framework: $Z_\phi: x \mapsto \mathbb{R}^b$

$\Phi(\theta)$ is the root of

$$\mathbb{E}_\theta[Z_\phi(x)]$$

i.e. $\Phi(\theta) \in \{\phi: \mathbb{E}_\theta[Z_\phi(x)] = 0\}$

Inference:

① Z_ϕ identifies $\Phi(\theta)$ if it has a unique root.

② $\phi_n \in \{\phi: \underbrace{P_n Z_\phi}_{\frac{1}{n} \sum_{i=1}^n Z_\phi(x_i)} = 0\}$

$\mathbb{E}_{x|1}[\text{Sample median}]$

If $\Phi(\theta) := \text{median}_\theta(X)$ is unique

then $0 = P_\theta[1\{X > \Phi(\theta)\}] - P_\theta[1\{X < \Phi(\theta)\}]$

$= P_\theta[\underbrace{\text{sign}(X - \Phi(\theta))}_{Z_\phi := \text{sign}(x - \phi)}]$

$Z_\phi := \text{sign}(x - \phi)$



$$\phi_n \in \{ \phi : P_n Z_\phi (= \frac{1}{n} \sum_{i=1}^n \text{sign}(X_i - \phi)) = 0 \}$$

matches the sample median

Ex. [Sample mean] $m\text{-est}: m_\phi = -|x - \phi|^2$

Ex. [sample median]

also an M-est: $m_\phi = -|x - \phi|$.

Ex. [Pearson method of moments]

Suppose $\Theta = \mathbb{R}^k$ and $\Phi: \Theta \rightarrow \Theta$, i.e.

we are interested in estimating θ_0 .

Suppose $X \sim P_{\theta_0}$ and the first k moments

$\begin{pmatrix} \mathbb{E}_\theta[X] \\ \vdots \\ \mathbb{E}_\theta[X^k] \end{pmatrix}$ are well-defined
and finite

$$\text{Then } P_{\theta_0} X (= \mathbb{E}_{\theta_0}[X]) \approx \frac{1}{n} \sum_{i=1}^n X_i$$

$$P_{\theta_0} X^2 \approx \frac{1}{n} \sum_{i=1}^n X_i^2$$

$$P_{\theta_0} X^k \approx \frac{1}{n} \sum_{i=1}^n X_i^k$$

Then the MoM est. is the root of

$$\left(\begin{array}{c} \frac{1}{n} \sum_{i=1}^n X_i - \mathbb{E}_\theta[X] \\ \frac{1}{n} \sum_{i=1}^n X_i^2 - \mathbb{E}_\theta[X^2] \\ \vdots \\ \frac{1}{n} \sum_{i=1}^n X_i^k - \mathbb{E}_\theta[X^k] \end{array} \right)$$

the blue parts
are known
functions of θ .

Then, Moll est. are Z_{est} :

$$Z_{\theta} : x \mapsto \begin{pmatrix} x - \mathbb{E}_{\theta}[x] \\ x^2 - \mathbb{E}_{\theta}[x^2] \\ \vdots \\ x^K - \mathbb{E}_{\theta}[x^K] \end{pmatrix}$$

plz. verify it

583

Ex. [M -est., but more general (U -stat)]

$$\text{PLM} : Y_i = X_i^T \beta_0 + g(U_i) + \varepsilon_i$$

outcome linear non-linear noise
 ↑ ↑ ↑ ↑
 unknown params.

covariates

a benchmark "semi-parametric" model since

$\{\beta_0 \in \mathbb{R}^d \text{ is finite-dim}$

$\{g \text{ is infinite-dim}$

$$\text{Goal of interest} : \Phi(\theta) = \beta_0$$

(β_0, g)

Difference-based estimator:

$$Y_i = X_i^T \beta_0 + g(U_i) + \varepsilon_i$$

$$Y_j = X_j^T \beta_0 + g(U_j) + \varepsilon_j$$

Pairwise diff

$$Y_i - Y_j = (X_i - X_j)^T \beta_0 + [g(U_i) - g(U_j)] + \varepsilon_i - \varepsilon_j$$

Then, if $U_i - U_j \approx 0$ and g is smooth, then

$$[g(U_i) - g(U_j)] \approx 0$$

$$Y_i - Y_j \approx (X_i - X_j)^T \beta_0 + \varepsilon_i - \varepsilon_j$$

$$\arg \min_{\beta \in \mathbb{R}^d} \left\{ \sum_{i < j} \underbrace{K_h(U_i - U_j)}_{\text{encouraging closeness of } U_i \text{ and } U_j} \underbrace{[Y_i - Y_j - (X_i - X_j)^T \beta]^2}_{\text{LSE loss}} \right\}$$

encouraging closeness of U_i and U_j

U -statistics for any β

$\{\beta : \sum_{i < j} K_h(\cdot) \dots\}$ is a U -process
indexed by β .

Full generalization of M- and Z-

M-est: $\hat{\theta}(\theta) \in \arg\max_{\theta \in S} M_\theta(\phi)$

population

special case: $P_\theta m_\phi$

$\hat{\theta}_n \in \arg\max_{\theta \in S} M_n(\phi)$

data-based

special case: $P_n m_\phi$

We believe as long as $M_n(\cdot)$ is close to $M_{\theta_0}(\cdot)$,
then their maximizer

$\hat{\theta}_n$ should be close to θ_0 ($\hat{\theta}(\theta_0)$)

- Z-est
- ① $\hat{\theta}_0 \in \{ \text{roots of } Z_{\theta_0}(\phi) \}$
 - ② $\hat{\theta}_n \in \{ \text{roots of } Z_n(\phi) \}$

Similarly, as long as $Z_{\theta_0}(\cdot)$ is close to
 $Z_n(\cdot)$, $\hat{\theta}_n$ should be close to θ_0 .

Remark:

(i) Usually, M-est. can be also formulated as an Z-est.

$$z_\theta: \phi \mapsto \nabla M_\theta(\phi)$$

$$z_n: \phi \mapsto \nabla M_n(\phi)$$

But not always true.

Manski's estimator:

Binary choice model goal of interest

$$Y_i = \mathbb{1}(X_i^\top \beta_0 + \varepsilon_i > 0)$$

↑ ↑ ↑
outcome covariates noise

$$\hat{\beta}_n := \underset{\beta \in \mathbb{R}^d}{\operatorname{argmax}} \sum_{i=1}^n \{ Y_i \mathbb{1}\{ X_i^\top \beta > 0 \} \}$$

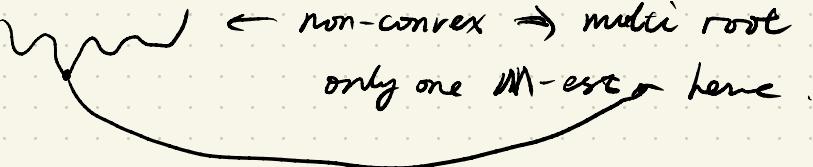
It doesn't have a natural Z-estimator.

(ii) Z-est. can always be formulated as
an M-est:

$$M_0 : \phi \mapsto \|Z_0(\phi)\|^2$$

$$M_n : \phi \mapsto \|Z_n(\phi)\|^2$$

Intuition: the global max is easier to be unique
than the root.

Ex.  ← non-convex \rightarrow multi root
only one M-est. here.

Q: "consistency"? "Rate of convergence"?
3. Dist.?

Is Q: Is the $\hat{\phi}_n$ above a consistent est. of ϕ_0 ??
i.e., $\hat{\phi}_n \xrightarrow{P} \phi_0$?

Example [1-dim Z-est.]

$\Phi : \mathbb{R} \mapsto \mathbb{R}$ s.t. $\phi_0 \in \mathbb{R}^1$.

$\phi_0 = \text{root of } Z_0(\phi)$ population score function

$\phi_n = \text{root of } Z_n(\phi)$ empirical score func.

Conditions:

(i) $\phi_0, \phi_n \in \underline{\mathbb{R}^1}$

(ii) [Point-wise consistency]: If given ϕ ,

$$Z_n(\phi) \xrightarrow{P} Z_0(\phi)$$

(iii) One of the following two is true

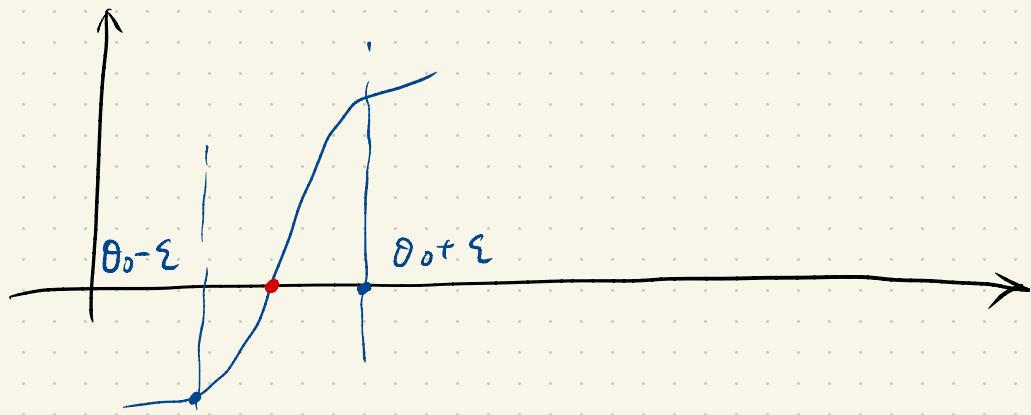
(a) Each $\phi \mapsto Z_n(\phi)$ is cont., and has exactly one root

or (b) $\phi \mapsto Z_n(\phi)$ is non-decreasing.

(iv) $\forall \varepsilon > 0, Z_0(\phi_0 - \varepsilon) < 0 < Z_0(\phi_0 + \varepsilon)$

Proof. We only show

(i)+(ii)+(iii) \Rightarrow (iv)



$$\begin{aligned}
 \text{Step 1 : } & P(Z_n(\phi_0 - \varepsilon) < 0, Z_n(\phi_0 + \varepsilon) > 0) \\
 & \leq P(\exists \text{ a root of } Z_n \text{ between } \phi_0 - \varepsilon \text{ and } \phi_0 + \varepsilon) \\
 & = P(\phi_n \in (\phi_0 - \varepsilon, \phi_0 + \varepsilon))
 \end{aligned}$$

On the other hand, we have known

$$(\text{iv}): Z_0(\phi_0 - \varepsilon) < 0, Z_0(\phi_0 + \varepsilon) > 0$$

$$(\text{ii}): Z_n(\phi_0 - \varepsilon) \xrightarrow{P} Z_0(\phi_0 - \varepsilon)$$

$$Z_n(\phi_0 + \varepsilon) \xrightarrow{P} Z_0(\phi_0 + \varepsilon)$$



$$P(Z_n(\phi_0 - \varepsilon) < 0, Z_n(\phi_0 + \varepsilon) > 0) \xrightarrow{P} 1$$

||

$$\phi_n \xrightarrow{P} \phi_0$$

Remark: This proof doesn't work for general \mathbb{R}^k param.

Real meat: $\phi_0 \in \operatorname{argmax} M_{\theta_0}(\phi)$

$\phi_n \in \operatorname{argmax} M_n(\phi)$ (*)

Goal: $\phi_n \xrightarrow{P} \phi_0$

Remark: (i) Normally, we can relax (*) to allow for computation error.

$$M_n(\phi_n) \geq \sup_{\phi} M_n(\phi) - o_p(1)$$

(ii) Identification: $\forall \varepsilon > 0$,

$$M_0(\phi_0) > \sup_{\phi: \|\phi - \phi_0\| > \varepsilon} M_0(\phi)$$

(iii) [Uniform consistency M_n to M_0]

$$\sup_{\phi \in S} |M_n(\phi) - M_0(\phi)| \xrightarrow{P} 0$$

Claim: (i) + (ii) + (iii) $\Rightarrow \phi_n \xrightarrow{P} \phi_0$

Map

$(\Phi, \Theta) \rightarrow$

$\Phi: \Theta \rightarrow \mathbb{R}^k$

$\underline{\Phi}(\theta)$

① maximizer

$M_\theta(\phi)$

$\underset{\phi \in S}{\operatorname{argmax}} P_\theta m_\phi$

② root of $P_\theta Z_\phi$

$Z_\theta(\phi)$

M - & Z - estimation

DGD:

$(P_0, \theta_0) \xrightarrow{\quad} \phi_0 = \underline{\Phi}(\theta_0)$

P_{θ_0}

Data

X_1, X_2, \dots, X_n
or
 P_n

est.
&
inference

$\phi_n \in$

$M_n(\phi)$

$\underset{\phi}{\operatorname{argmax}} P_n m_\phi$

$Z_n(\phi)$

root of $P_n Z_\phi$

Consistency & rate of convergence

$$(i) M_n(\phi_n) \geq \sup_{\phi} M_n(\phi) - \underline{O_p(1)}$$

minor computation error

(ii) $\forall \varepsilon > 0,$

$$M_0(\phi_0) > \sup_{\phi: \| \phi - \phi_0 \| > \varepsilon} M_0(\phi)$$

(iii) Uniform convergence

$$\sup_{\phi \in S} |M_n(\phi) - M_0(\phi)| \xrightarrow{P} 0$$

Goal: $\phi_n \xrightarrow{P} \phi_0.$

Proof: $0 \leq M_0(\phi_0) - M_0(\phi_n)$

$$\leq [M_0(\phi_0) - M_0(\phi_n)]$$

$$- [M_n(\phi_0) - \sup_{\phi \in S} M_n(\phi)] \stackrel{\leq 0}{\leftarrow}$$

$$\leq [M_0(\phi_0) - M_0(\phi_n)]$$

$$- [M_n(\phi_0) - M_n(\phi_n)] + O_p(1) \stackrel{(i)}{\nwarrow}$$

$$\begin{aligned}
&= \underbrace{[M_0(\phi_0) - M_n(\phi_0)]}_{\leq \sup_{\phi} |(M_n - M_0)(\phi)|} + \underbrace{[M_n(\phi_n) - M_0(\phi_n)]}_{\leq \sup_{\phi} |(M_n - M_0)(\phi)|} + o_p(1) \\
&\leq 2 \sup_{\phi} |(M_n - M_0)(\phi)| + o_p(1) \stackrel{(iii)}{=} o_p(1)
\end{aligned}$$

$$\Rightarrow M_0(\phi_0) - M_0(\phi_n) = o_p(1) \quad (*)$$

It remains to show

$$(*) \xrightarrow{(ii)} \phi_n \xrightarrow{P} \phi_0$$

Fix $\varepsilon > 0$, and let

$$S = M_0(\phi_0) - \sup_{\|\phi - \phi_0\| > \varepsilon} M_0(\phi) \stackrel{(ii)}{>} 0$$

Notice

$$\{\|\phi_n - \phi_0\| > \varepsilon\} \subset \{M_0(\phi_0) - M_0(\phi_n) \geq S\}$$



$$P_0(\|\phi_n - \phi_0\| > \varepsilon) \leq P_0(M_0(\phi_0) - M_0(\phi_n) \geq S)$$

$$\xrightarrow{n \rightarrow \infty} = 0$$

$$\phi_n \xrightarrow{P} \phi_0$$

Uniform consistency

$$\text{Setup: } M_0(\phi) = P_0 m_\phi$$

$$M_n(\phi) = P_n m_\phi = \frac{1}{n} \sum_{i=1}^n m_\phi(X_i)$$

Prop. As long as

$$(i) \forall \phi \in S, |P_0 m_\phi| < \infty$$

$$(ii) |S| < \infty$$

We have

$$\sup_{\phi \in S} |(P_n - P_0) m_\phi| \xrightarrow{P} 0$$

Proof. SLLN $\stackrel{(i)}{\downarrow}$ + (ii)

Union bound: as long as

$$A_1, A_2, \dots \subset \mathcal{X}$$

$$\text{then } \mu\left(\bigcup_{i=1}^{\infty} A_i\right) \leq \sum_{i=1}^{\infty} \mu(A_i)$$

Step 1:

By LLN, $\forall \varepsilon > 0$,

$$\begin{aligned} & P_0(|P_n - P_0| m_\phi > \varepsilon) \\ &= P_0\left[\left|\frac{1}{n} \sum_{i=1}^n (m_\phi(X_i) - P_0 m_\phi)\right| > \varepsilon\right] \end{aligned}$$

$\xrightarrow{\text{LLN}(i)}$ $\rightarrow 0$, as $n \rightarrow \infty$

$$\text{Step 2. } P_0 \left(\sup_{\phi \in S} |P_n - P_0| m_\phi > \varepsilon \right)$$

$$= P_0 \left(\bigcup_{\phi \in S} \{|P_n - P_0| m_\phi > \varepsilon\} \right)$$

$$\leq \sum_{\phi \in S} P_0(|P_n - P_0| m_\phi > \varepsilon)$$

union
bound

$$\leq |S| \cdot \sup_{\phi \in S} P_0(|P_n - P_0| m_\phi > \varepsilon)$$

$$\xrightarrow{n \rightarrow \infty} 0$$

$$\downarrow n \rightarrow \infty$$

0

For studying

$$\sup_{\phi \in S} |P_n - P_0| m_\phi$$

it suffices to study

$$\{m_\phi : \phi \in S\}$$

Def. [P₀-Glivenko-Cantelli]

The class of functions

$$\{m_\phi : \phi \in S\}$$

is called a P₀-GC class/set

if it satisfies

$$\sup_{\phi \in S} |P_n - P_0| m_\phi = o_p(1)$$

Remark G & C independently studied

$$\sup_{t \in \mathbb{R}} |\underbrace{F_n(t)}_{\text{empirical COF}} - \underbrace{F_0(t)}_{\text{population COF}}| \xrightarrow{P} 0$$

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}\{X_i \leq t\} \xrightarrow{\text{P}} P(X \leq t)$$

$$\underbrace{\frac{1}{n} \sum_{i=1}^n \mathbf{1}\{X_i \leq t\}}_{\text{"}} \xrightarrow{\text{P}} P(X \leq t)$$

$$P_n \mathbf{1}(X \leq \cdot) \xrightarrow{\text{P}} P_0 \mathbf{1}(X \leq \cdot)$$

$$\Leftrightarrow \sup_{\phi \in \mathbb{R}} |P_0 - P_0| m_\phi \xrightarrow{P} 0$$

$$\text{w.r.t. } \{m_\phi := \mathbf{1}\{X \leq \phi\}, \phi \in \mathbb{R}\}$$

Remark. Methods to prove a certain function class is P₀-GC:

(i) Symmetrization + VC-type bound

(ii) Martingale theory [520s]

(iii) Entropy argument [581]

metric-entropy, bracketing-entropy --

Def: [bracketing entropy]

Setup: $\mathcal{F} = \{f \in \mathcal{G}\}$

$f: X \rightarrow \mathbb{R}^1$ real line

(i) [Bracket] Given any two functions

l and u [$l < u$]

in $L^1(P_0)$ [i.e., $\int |l(x)| P_0(dx) < \infty$

$\int |u(x)| P_0(dx) < \infty$]

the bracket of l and u ,

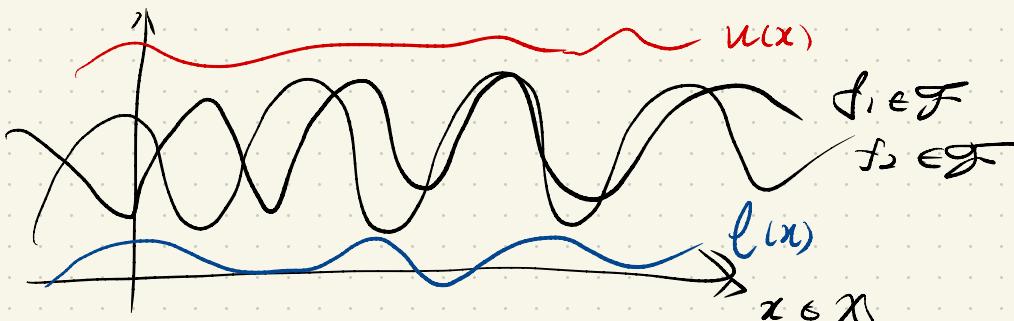
$[l, u]$

containing all functions

$f \in \mathcal{F}$

s.t.

$l(x) \leq f(x) \leq u(x), \forall x \in X$.



(ii) A bracket $[l, u]$ is said to be an ϵ -bracket if

$$\|u - l\|_{L^1(P_0)} \leq \epsilon.$$

$$\|f\|_{L_1(P_0)} := \int |f(x)| P_0(dx)$$

(iii) Bracketing entropy of \mathcal{F} .

$$N_{[\cdot]}(\varepsilon, L^1(P_0), \mathcal{F})$$

represents the smallest number
of ε -brackets needed to cover \mathcal{F} .

Remark. Notice the l 's and n 's do not have to be in \mathcal{F} , but they should $\in L^1(P_0)$.

Thm. [P_0 -GC using bracketing entropy]

Suppose, $\forall \varepsilon > 0$,

$$N_{[\cdot]}(\varepsilon, L^1(P_0), \mathcal{F}) < \infty$$

Then, \mathcal{F} is P_0 -GC, i.e.,

$$\begin{aligned} \|P_n - P_0\|_{\mathcal{F}} &= : \sup_{f \in \mathcal{F}} |P_n f - P_0 f| \\ &= o_p(1) \end{aligned}$$

Proof. (i) with each ε -bracket,

f 's are nearly the same (identical)

(ii) \exists only finitely many ε -brackets

(iii) Union bound proof.

Fix $\varepsilon > 0$. By cond.,

$$\exists [l_j, u_j] \cdot j=1, 2, \dots, N \boxed{[l_j, u_j]} \leftarrow N$$

s.t. it covers \mathcal{F} .

Define

$$A_{n,\varepsilon} := \left\{ \sup_{f \in \mathcal{F}} (P_n - P)f > 2\varepsilon \right\}$$

$$B_{n,\varepsilon} := \left\{ \inf_{f \in \mathcal{F}} (P_n - P)f < -2\varepsilon \right\}$$

By symmetry, only study $A_{n,\varepsilon}$.

Fix $f \in \mathcal{F}$, By cond.,

\exists one $j \in \{1, 2, \dots, N\}$ s.t.

$$l_j \leq f \leq u_j$$

implying

$$P_n f \leq P_n u_j$$

$$P_n f \geq P_n u_j - P_n(u_j - l_j)$$

$$(P_n - P_0) f \leq P_n u_j - P_n u_j + P_n(u_j - l_j)$$

$$= (P_n - P_0) u_j + P_0(u_j - l_j)$$

$$= (P_n - P_0) u_j + \underbrace{P_0 |u_j - l_j|}_{\parallel u_j - l_j \parallel_{L^2(P_0)}}$$

$$= \parallel u_j - l_j \parallel_{L^2(P_0)}$$

$$\leq \varepsilon$$

$$\leq (P_n - P_0) u_j + \varepsilon$$

$$\begin{aligned}
 & P_0(A_{n,\varepsilon}) \\
 &= P_0\left(\sup_{f \in \mathcal{F}} (P_n - P_0)f > 2\varepsilon\right) \\
 &\leq P_0\left(\sup_{j \in \{1, 2, \dots, n\}} (P_n - P_0)u_j > \varepsilon\right) \\
 &\leq \sum_{j=1}^n P_0((P_n - P_0)u_j > \varepsilon) = o(1)
 \end{aligned}$$

□

Real Q: How to show a class of functions have finite bracketing entropy

$N_{[]}(\cdot) ???$

Thm [Example 19.8 Van der Vaart]

① Suppose

$$\mathcal{F} := \{f_\phi, \phi \in K \subset \mathbb{R}^d\}$$

where K is compact.

② \forall given x ,

$\phi \mapsto f_\phi(x)$ is continuous w.r.t. ϕ .

③ [envelope] Suppose

\exists a function F satisfying

$$(a) \sup_{\phi \in K} |f_\phi(x)| \leq F(x), \quad \forall x$$

$$(b) P_0|F| = P_0 F < \infty$$

Then $\forall \varepsilon > 0$, $N_{[]}(\varepsilon, L^2(P_0), \mathcal{F})$ is finite.

Proof HW!

Ex. [logistic reg.]

Setup: $X = (Z, Y)$

Data covariates outcome

GLM: $Y|Z \sim \text{Ber}(g_{\Phi(\theta)}(Z))$

$$\text{w. } g_\phi(Z) = \frac{1}{1 + \exp(-\phi^T Z)}$$

$$Z \sim P_\theta \left\{ \begin{array}{l} \text{Gaussian} \\ \text{others} \end{array} \right.$$

n IID obs. of (Z, Y)

X_1, X_2, \dots, X_n

Goal: Estimate ϕ

Assumption: $\Phi(\theta_0) = \phi_0 \in K \subseteq \mathbb{R}^d$ compact

Procedure: MLE / M-est.

$$\hat{\phi}_n = \underset{\phi \in K}{\operatorname{argmax}} P_n m_\phi$$

$$\text{with } m_\phi: z = (Z, Y) \mapsto y \log g_\phi(z) + (1-y) \log(1-g_\phi(z))$$

$$= y \log(g_\phi(z)) + \log(1-g_\phi(z))$$

$$= y \cdot \phi^T z - \log(1 + \exp(\phi^T z))$$

"

① $\phi \in K$ compact ✓

② $\forall x = (y, z)$, it is cont. w.r.t. ϕ . ✓

③ It remains to find the envelope F :

$$\sup_{\phi \in K} |m_\phi(x)| \leq |y \cdot \phi^T z| + |\log(1 + \exp(\phi^T z))|$$

$$\leq |\phi^T z| + \log(1 + \exp(|\phi^T z|))$$

$$\begin{cases} 0 \leq \log(1 + \exp(t)) \\ \leq \max(t, 0) + 1 \\ \forall t \end{cases}$$

$$\leq |\phi^T z| + \max(|\phi^T z|, 0) + 1$$

$$= 2|\phi^T z| + 1$$

$$[\text{Cauchy-Schwarz}] \leq 2\|\phi\| \|z\| + 1$$

$$\sup_{\phi \in K}$$

since K compact,
it is finite, $< \infty$.

$$= 2\|z\| \sup_{\phi \in K} \|\phi\| + 1$$

$$:= F, : x \mapsto 2\|z\| \sup_{\phi \in K} \|\phi\| + 1$$

If $\|F\|_{L^1(P_0)} < \infty$, then done.

$$\Leftrightarrow P_0\|z\| < \infty$$

$$\Leftrightarrow P_0 z_j^2 < \infty, \forall j = 1, \dots, d.$$

By Vdr

$$N_{\mathcal{D}}(\varepsilon, L'(P_0), \Sigma_{\mathbf{m}\phi: \phi \in K}) < \alpha$$

ULLN is true

$P_0 m_\phi$ is convex \leftarrow verify

$$\phi_n \xrightarrow{P} \phi_0$$

3.4 ASN of M- & Z- est.

$$\sqrt{n}(\phi_n - \phi_0) \Rightarrow N(0, \sigma^2)$$

asymptotic normality, a.k.a. ASN

To establish ASN of ϕ_n , usually we aim to show

Asymptotic linear estimator (ALE)

$$\sqrt{n}(\phi_n - \phi_0) = \boxed{\frac{1}{\sqrt{n}} \sum_{i=1}^n z_i} + \boxed{R_n} + o_p(1)$$

\uparrow
linear expansion
/ Hajak projection

Donsker property

$$\sqrt{n}(P_n - P)(\hat{f}_n - \hat{g}_n)$$

$$\text{with } d(\hat{f}_n, \hat{g}_n) \xrightarrow{n \rightarrow \infty} 0$$

The most challenging part is to show

$$R_n = o_p(1) \quad (\text{Take supremum})$$

stochastic equicontinuity

$$\lim_{\varepsilon \rightarrow 0} \sup_{f, g : d(f, g) < \varepsilon} \sqrt{n} (P_n - P)(f - g) = 0$$

$\exists f$ s.t. $\boxed{\quad}$ is true is called
Po-Donsker

Example (mean absolute deviation, David Pollard)

$$u_n := \frac{1}{n} \sum_{i=1}^n |X_i - \bar{X}_n| \quad \text{Robust statistics}$$

[Setup: $X_1, \dots, X_n \stackrel{iid}{\sim} P_0$,
 $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$]

1st [Consistency]

$$\mu_n \xrightarrow{P} \mu_0 := E[|X - E[X]|]$$

2nd Mean absolute deviation is ASN

Goal: $\sqrt{n}(\mu_n - \mu_0) \Rightarrow N(0, \sigma^2)$

Proof: Def. $\mathcal{F} := \{ \underbrace{|x-t|}_{f_t(x)} : t \in \mathbb{R} \}$

$$\mu_n = P_n f_{\bar{X}_n}$$

$$\mu_0 = P_0 f_{\mu_0} (= E[X])$$

$$\sqrt{n}(\mu_n - \mu_0) = \sqrt{n}(P_n f_{\bar{X}_n} - P_0 f_\mu)$$

$$= \boxed{\sqrt{n}(P_n - P_0)f_\mu} \leftarrow \text{CLT}$$

$$+ \sqrt{n}(P_n f_{\bar{X}_n} - P_n f_\mu)$$

$$= \boxed{\sqrt{n}(P_n - P_0)f_\mu} \leftarrow \text{CLT} + \boxed{\sqrt{n}(\psi(\bar{X}_n) - \psi(\mu))} \stackrel{\psi: \text{smooth}}{\approx} \text{Poft} \leftarrow \text{Delta method}$$

$$+ \boxed{\sqrt{n}(P_n - P_0)(f_{\bar{X}_n} - f_\mu)}$$

↑
Donsker

$$\Rightarrow N(0, \sigma^2)$$

$$= \sqrt{n}(P_n - P_0)(f_\mu(x) + \boxed{\psi'(\mu) \cdot x}) \leftarrow \begin{array}{l} \text{Delta method main} \\ \text{term by Taylor} \\ \text{expansion} \end{array}$$

$$+ \boxed{O_p(1)}$$

↓
Delta method

$$+ \boxed{O_p(1)} \quad \text{STAT 582}$$

↑
Donsker's $O_p(1)$

$$\Rightarrow N(0, \boxed{\sigma^2}) \quad \text{Please calculate it}$$

ASV of Z-est.

Cramer 1930

Setup [simple]

Cond. too strong
and eliminate very important
application like quantile regression.

$$\phi_0 \in \mathbb{R}$$

$$\phi_0: Z_0(\phi_0) = 0$$

$$\phi_n: Z_n(\phi_n) = 0$$

$$\text{Goal: } \sqrt{n}(\phi_n - \phi_0) \Rightarrow N(0, \sigma^2)$$

513-level:

① Taylor expansion on Z_n

smoothness
is required
naguthy
might not be 0

② Cramer-type cond.

③ Smoothness on Z_n 's function

They are strong conditions and not necessary.

Claim: ① - ③ cannot handle many estimators, such as quantile reg:

$$\underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{i=1}^n |Y_i - X_i^\top \beta|$$

Insight: too much assumption on Z_n which can be very non-smooth.

Instead of Z_n , we can look at Z_0 .

Because Z_0 we can add expectation,

which will smoothen the original function.

Use Donsker to connect Z_n and Z_0 .

Details (heuristic)

$$O = -Z_0(\phi_0)$$

$$= [Z_n(\phi_0) - Z_0(\phi_0)] - Z_n(\phi_0)$$

$$= [Z_n(\phi_0) - Z_0(\phi_0)]$$

$$\begin{aligned} Z_n(\phi_n) &= 0 \\ &\quad + [Z_n(\phi_n) - Z_n(\phi_0)] \end{aligned}$$

$$I = [Z_n(\phi_0) - Z_n(\phi_0)] \rightarrow CLT$$

$$II = Z_0(\phi_n) - Z_0(\phi_0) \rightarrow \begin{array}{l} \text{usually smooth} \\ \rightarrow \text{Taylor expansion} \end{array}$$

$$III = \boxed{\begin{aligned} &+ Z_n(\phi_n) - Z_n(\phi_0) \\ &- Z_0(\phi_n) + Z_0(\phi_0) \end{aligned}} \rightarrow \text{a Donsker term}$$

$$\text{When } Z_n(\phi) = P_n Z_\phi$$

$$Z_0(\phi) = P_0 Z_\phi$$

the last term = $(P_n - P_0)(Z_{\phi_n} - Z_{\phi_0})$

Step I, Term I = $(P_n - P_0) Z_{\phi_0}$

$\Rightarrow \sqrt{n} \text{Term I} \Rightarrow N(0, \text{Var}(Z_{\phi_0}))$

Step II, Term II = $Z_0(\phi_n) - Z_0(\phi_0)$

$\Rightarrow \text{Term II} = (\phi_n - \phi_0) \dot{Z}_0(\phi_0)$
 $+ \frac{1}{2} (\phi_n - \phi_0)^2 \ddot{Z}_0(\tilde{\phi}_n)$

for some $\tilde{\phi}_n$ between ϕ_0 and ϕ_n .

$$= (\phi_n - \phi_0) \dot{Z}_0(\phi_0) + o_p(\phi_n - \phi_0)$$

If ① $\phi_n \xrightarrow{P} \phi_0$

② $\sup_{\phi \in U(\phi_0)} |\ddot{Z}_0(\phi)| < \infty$
small neighbourhood

Step 3, We need to show

$$\sqrt{n} \cdot \text{Term III} = o_p(1)$$

An incorrect heuristic to step 3.

LéCam's thinking: instead of having

a random seq. of $\{\phi_n\}$, let's think about
" $\{\phi_n\}$ to be deterministic"

INCORRECT but nearly the sample splitting.

Under the Incorrect assumption:

$$\text{Term II} = (P_n - P_0) (\bar{Z}_{\phi_n} - \bar{Z}_{\phi_0})$$

Fixing any $t > 0$,

chebyshev $P_0 [\sqrt{n} \text{Term II} > t]$

$$\leq \frac{\text{Var}(\sqrt{n}(P_n - P_0)(\bar{Z}_{\phi_n} - \bar{Z}_{\phi_0}))}{t^2}$$

$$= \frac{\text{Var}(\bar{Z}_{\phi_n}(X_1) - \bar{Z}_{\phi_0}(X_1))}{t^2}$$

please verify

Then, under certain conditions,

$$\text{Var}(\bar{Z}_{\phi_n}(\cdot) - \bar{Z}_{\phi_0}(\cdot)) \rightarrow 0$$

$$\text{as } \phi_n \rightarrow 0$$

e.g. Lipschitz cond:

$$|\bar{Z}_\phi(x) - \bar{Z}_{\phi_0}(x)| \leq \|\phi - \phi_0\| \cdot G(x)$$

$$\text{and } P_0 G^2 < \infty$$

Step 4.

$$0 = \text{Term I} + \text{Term II} + \text{Term III}$$

$$\Rightarrow \phi_n - \phi_0 = \frac{-(P_n - P_0) \bar{Z}_{\phi_0}}{\dot{Z}_0(\phi_0) + o_p(1)} \leftarrow \text{Term I}$$

$$+ o_p(n^{-1/2}) \leftarrow \text{Term II}$$

$$\Rightarrow \sqrt{n}(\hat{\phi}_n - \phi_0) = \frac{-\sqrt{n}(P_n - P_0) Z_{\phi_0}}{\dot{Z}_{\phi_0}(\phi_0)} \xrightarrow{ASL} + o_p(1)$$

$$\Rightarrow N(0, \boxed{\frac{\text{var}_0(Z_{\phi_0})}{(\dot{Z}_{\phi_0}(\phi_0))^2}})$$

$$\ln MLE, = \frac{1}{I(\theta_0)}$$

The following 2 Thms will not be proved.

Thm [General Z-est. ASN]

① $Z_0 = P_0 Z_\phi$, $Z_n = P_n Z_\phi$, and $\forall \phi$ in an open
subset of \mathbb{R}^d

$$Z_\phi: X \mapsto \mathbb{R}^d$$

$$② E_0 \|Z_{\phi_0}(X)\|^2 < \infty$$

③ $\phi \mapsto P_0 Z_\phi$ is differentiable (Taylor expansion)
at ϕ_0 with nonsingular
Jacobian V_{ϕ_0} . \longrightarrow ln MLE, Fisher info matrix

④ Assume $\exists G: X \mapsto \mathbb{R}$ (smoothness for Donsker)
s.t. ① $P_0 G^2 < \infty$

② $\forall x \in X, \forall \phi, \tilde{\phi} \in U(\phi_0),$

$$\|Z_\phi(x) - Z_{\tilde{\phi}}(x)\| \leq \|\phi - \tilde{\phi}\| \cdot G(x)$$

⑤ $\{\phi_n\}$ is a sequence of est. ϕ_0 s.t.

$$P_n Z_{\phi_n} = o_p(n^{-1/2})$$

and $\phi_n \xrightarrow{P} \phi_0$.

Claim: Under ①-⑤,

$$\sqrt{n}(\phi_n - \phi_0) \xrightarrow{} N(0, V_{\phi_0}^{-1} \cdot P_0 [Z_{\phi_0} Z_{\phi_0}^T] [V_{\phi_0}^{-1}])$$

Thm [ASN of M-est.] Assume

① $\phi_0 = \operatorname{argmax} P_0 m_\phi$

$[\phi_n \approx \operatorname{argmax} P_n m_\phi]$

$m_\phi: \mathcal{X} \mapsto \mathbb{R}$

② $\forall \phi \in \text{open subset } \subset \mathbb{R}^d$,

$\phi \mapsto m_\phi(x)$ is differentiable at ϕ_0 for
 P_0 -almost everywhere x .

Note: think about $|x|$ has measure 0
∴ still satisfies cond. ②

We take

\dot{m}_{ϕ_0} to be the derivative.

③ $\forall \phi, \tilde{\phi} \in U(\phi_0)$, assume

$\exists G: \mathcal{X} \mapsto \mathbb{R}$ s.t.

① $P_0 G^2 < \infty$

② $|m_\phi(x) - m_{\tilde{\phi}}(x)| \leq \|\phi - \tilde{\phi}\| G(x), \forall x \in \mathcal{X}$.

Note: even weaker than Z because M-est. acts like integral of Z-est.

④ Assume \exists a non-singular symmetric matrix

V_{ϕ_0} s.t.

$$\lim_{\varepsilon \rightarrow 0} \frac{\sup_{\|h\|=1} |P_0 m_{\phi_0 + \varepsilon h} - P_0 m_{\phi_0} - \frac{1}{2} \varepsilon^2 h^\top V_{\phi_0} h|}{\varepsilon^2} = 0$$

[We will use QMD to verify]

⑤ $P_n m_{\phi_n} \geq \sup_{\phi} P_n m_{\phi} - o_p(n^{-1})$
& $\phi_n \xrightarrow{P} \phi_0$

Claim: ① - ⑤ imply

$$\sqrt{n}(\phi_n - \phi_0) \xrightarrow{\text{ALE}} N(0, V_{\phi_0}^{-1} [P_0 m_{\phi_0}] [V_{\phi_0}^{-1}]^\top)$$

$$\left(\begin{array}{l} \\ \end{array} \right) = V_{\phi_0}^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n m_{\phi_0}(x_i) + o_p(1)$$

Slick by

As long as it is not ASL, Bootstrap cannot work.

Bootstrap consistency \approx ASL.

BS can improve finite sample accuracy.

BS can accelerate rate.

Chapter 3.5 Parametric models

(MLE + QMD)



K-L divergence Hellinger distance

Setup: $\mathcal{P} := \{ P_\theta : \theta \in \Theta \subset \mathbb{R}^d \xrightarrow{\text{open}} \xrightarrow{\text{finite dim}}$

MLE, maximize log-likelihood.

is a M-est:

$$M_{\theta_0} : \Theta \mapsto E_\theta \left[\log \frac{dP_\theta}{d\mu}(x) \right]$$

① μ is the reference measure { Lebesgue
Counting }

② Choose μ to be Lebesgue measure.

$\frac{dP_\theta}{d\mu}$ is the pdf

③ Choose μ to be counting measure,

$\frac{dP_\theta}{d\mu}$ is the pmf.

We'll write

$$P_\theta = \frac{dP_\theta}{d\mu} \quad (\text{w.r.t. } \mu)$$

M_{θ_0} will give us Z_{θ_0}

$$\begin{aligned} Z_{\theta_0} &: \Theta \mapsto \nabla_{\theta} M_{\theta_0} \\ &= \nabla_{\theta} E_\theta [\log P_\theta(x)] \end{aligned}$$

High level motivation of QMD

∇_{θ} and \int

Interchange the position

$$\text{Ex} \quad Z_{\theta_0}(\theta) = \mathbb{E}_{\theta_0}[Z_{\theta}(x)]$$

$$\begin{aligned} Z_{\theta} := x &\mapsto \nabla_{\theta} \log P_{\theta}(x) \\ &= \boxed{\frac{\dot{P}_{\theta}(x)}{P_{\theta}(x)}} l_{\theta}(x) \\ &(\dot{P}_{\theta}(x) = \nabla_{\theta} P_{\theta}(x)) \end{aligned}$$

Claim [513] "In general,"

$$Z_{\theta_0}(\theta_0) = 0$$

Proof. $Z_{\theta_0}(\theta_0) = \mathbb{E}_{\theta_0}[l_{\theta_0}(x)]$

$$= \int l_{\theta_0}(x) P_{\theta_0}(x) \mu(dx)$$

$$= \int \dot{P}_{\theta_0}(x) \mu(dx)$$

cond
is very
naughty \rightsquigarrow $= \int \nabla_{\theta} P_{\theta_0}(x) \mu(dx)$

$$\nabla_{\theta} \underbrace{\int P_{\theta_0}(x) \mu(dx)}_{Y} = 0$$

Claim 2 [STAT 513 level]

ASN of MLE :

We can formulate the problem as

Z-est:

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \Rightarrow N\left(0, \frac{P_{\theta_0}[\dot{Z}_{\theta_0}]^2}{[\dot{Z}_{\theta_0}(\theta_0)]^2}\right)$$

[when $\theta_0 \in \mathbb{R}$]

Intuitively,

$$\frac{P_{\theta_0}[\dot{Z}_{\theta_0}]^2}{[\dot{Z}_{\theta_0}(\theta_0)]^2} = \frac{1}{I(\theta_0)} \quad (*)$$

← FIN at θ_0

To prove (*), 513 tells us

$$\dot{Z}_{\theta_0}(\theta_0) = \frac{d}{d\theta} \int \dot{l}_{\theta_0}(x) P_{\theta_0}(dx)$$

$$= \int \frac{d}{d\theta} \dot{l}_{\theta_0}(x) P_{\theta_0}(dx)$$

$$\begin{aligned} \text{Interchange} &= \int \frac{\ddot{P}_{\theta_0}(x) P_{\theta_0}(x) - \dot{P}_{\theta_0}^2(x)}{[P_{\theta_0}(x)]^2} P_{\theta_0}(dx) \\ &\qquad\qquad\qquad P_{\theta_0}(x) \mu(dx) \end{aligned}$$

$$= \int \ddot{P}_{\theta_0}(x) \mu(dx)$$

$$\begin{aligned} \frac{d^2}{d\theta^2} \int_{\theta_0} P_{\theta_0}(x) \mu(dx) &= - \boxed{\int [\dot{l}_{\theta_0}(x)]^2 P_{\theta_0}(dx)} \\ &\qquad\qquad\qquad \text{← var of score function} \\ &\qquad\qquad\qquad "FIN" \end{aligned}$$

$$= -FIN$$

Chapter 3.6 QMD

Quadratic Mean Differentiability

Note: $\nabla_{\theta} (\sqrt{P_{\theta}(x)}) = \frac{1}{2} \cdot \frac{\dot{P}_{\theta}(x)}{\sqrt{P_{\theta}(x)}}$

$$= \frac{1}{2} \cdot \frac{\dot{P}_{\theta}(x)}{P_{\theta}(x)} \cdot \sqrt{P_{\theta}(x)}$$

$$= \frac{1}{2} \cdot \dot{e}_{\theta}(x) \cdot \sqrt{P_{\theta}(x)}$$

Def [QMD]

(i) The root density

$$\theta \mapsto \sqrt{P_{\theta}}$$

is called QMD at θ if

\exists a function e_{θ}

perturbation

s.t.

$$\lim_{\varepsilon \rightarrow 0} \sup_{\|h\|=1} \int \left[\frac{\sqrt{P_{\theta+h}(x)} - \sqrt{P_{\theta}(x)}}{\varepsilon} \right]^2 \mu(dx) = 0$$

mean differentiability

quadratic

related to Hellinger distance
(metric)

$$\forall p, q < \mu, \int (\sqrt{p} - \sqrt{q})^2 \mu(dx)$$

$$p(x) = \frac{dp}{d\mu}, q = \frac{dq}{d\mu}$$

(ii) A model $\{P_{\theta} : \theta \in \Theta\}$ is QMD at θ
if $\theta \mapsto \sqrt{P_{\theta}}$ is QMD at θ .

This model is QMD if it is QMD at $\forall \theta \in \Theta$.

Remark: An equivalent expression for QMD at θ is

$$\int \left[\sqrt{P_{\theta+h}} - \sqrt{P_\theta} - \frac{1}{2} h^T \dot{\ell}_\theta \cdot \sqrt{P_\theta} \right]^2 \mu(dx)$$

$$= o(\|h\|^2) \text{ as } h \rightarrow 0 \quad [\text{VdV book}]$$

uniform [Please verify]

Thm. [1st part of VdV's Thm. 7.2]

Suppose \mathbb{D} is open subset of \mathbb{R}^d and

$\dot{\ell}_\theta(\cdot)$ is
the function
is QMD

$\{P_\theta : \theta \in \mathbb{D}\}$ is QMD at θ . Then

$$\textcircled{1} \quad P_\theta \dot{\ell}_\theta = 0 \quad (\text{1st Eq})$$

$$\textcircled{2} \quad I_\theta = P_\theta [\dot{\ell}_\theta \dot{\ell}_\theta^T] \text{ exists.}$$

Proof of $\textcircled{1}$ QMD at θ gives us

$\textcircled{1} \quad \frac{\sqrt{P_{\theta+\epsilon h}(x)} - \sqrt{P_\theta(x)}}{\epsilon} \text{ converge in QM}$

$$\text{to } \frac{1}{2} h^T \dot{\ell}_\theta \cdot \sqrt{P_\theta}$$

$\textcircled{2} \quad \sqrt{P_{\theta+\epsilon h}(x)} - \sqrt{P_\theta(x)} \text{ converges in QM}$
to 0.

To prove $P_\theta \dot{\ell}_\theta = 0$, it suffices to show

$$\underline{P_\theta[h^T \dot{\ell}_\theta]} = 0, \forall h \text{ s.t. } \|h\|=1$$

$$= \int h^T \dot{\ell}_\theta P_\theta d\mu$$

$$= \int [h^T \dot{\ell}_\theta \cdot \sqrt{P_\theta}] \sqrt{P_\theta} d\mu$$

$$\begin{aligned}
 &= \lim_{\varepsilon \rightarrow 0} 2 \int \left[\frac{\sqrt{P_{\theta+\varepsilon h}} - \sqrt{P_\theta}}{\varepsilon} \right] \frac{\sqrt{P_{\theta+\varepsilon h}} + \sqrt{P_\theta}}{2} d\mu \\
 \text{HW} \quad &= \lim_{\varepsilon \rightarrow 0} \int \frac{P_{\theta+\varepsilon h} - P_\theta}{\varepsilon} d\mu \\
 &= 0 \quad \square
 \end{aligned}$$

Thm [ASIN of MLE under QMD]

Suppose

- ① $\{P_\theta : \theta \in \Theta\}$ is QMD at θ_0 , which is an inner point of Θ .
- ② \exists a function G satisfying
 - (a) $P_0 G^2 < \infty$
 - (b) $\forall \theta_1, \theta_2 \in U(\theta_0)$, $\forall x$

$$|\log P_{\theta_1}(x) - \log P_{\theta_2}(x)| \leq G(x) \|\theta_1 - \theta_2\|$$
- ③ I_{θ_0} is non-singular (invertible) and the MLE $\hat{\theta}_n$ is consistent.

Then, we have asymptotic normality.

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \Rightarrow N(0, I(\theta_0)^{-1})$$

Proof. We only prove ④ is M-est. Thm,

$$\sup_{\|h\|=1} \left| P_0 M_{\theta_0+\varepsilon h} - P_0 M_{\theta_0} + \frac{1}{2} \varepsilon^2 h^T V_{\theta_0} h \right| = o(\varepsilon^2)$$

as $\varepsilon \rightarrow 0$.

$$\begin{aligned}
 \square &= P_0 \left[\log \underbrace{P_{0+\varepsilon h}}_{\overset{\leftarrow}{P_\varepsilon}} - \log \underbrace{P_0}_{\overset{\leftarrow}{P_0}} \right] \\
 &= P_0 \left[\log P_\varepsilon - \log P_0 \right] \\
 &= P_0 \left[2 \log \sqrt{P_\varepsilon} - 2 \log \sqrt{P_0} \right] \\
 &= 2 P_0 \left[\log \sqrt{P_\varepsilon} - \log \sqrt{P_0} \right] \\
 &= 2 P_0 \left[\log \frac{\sqrt{P_\varepsilon}}{\sqrt{P_0}} - \log 1 \right] \\
 &= 2 \int [f(W_\varepsilon) - f(0)] P_0(dx)
 \end{aligned}$$

where $f(w) = \log(1+w)$

$$W_\varepsilon = \sqrt{P_\varepsilon/P_0} - 1$$

Now do Taylor expansion on $f(\cdot)$: Lyapunov version of Taylor expansion

$$f(w) - f(0) = w - \frac{w^2}{2} + \underline{w^2 \cdot r(w)}$$

with $r(w) \rightarrow 0$ as $w \rightarrow 0$ "o(w^2)

$$\begin{aligned}
 \Rightarrow \square &= \boxed{2 \int W_\varepsilon \cdot P_0(dx)} = \int (\sqrt{\frac{P_\varepsilon}{P_0}} - 1) P_0 \mu(dx) \\
 &\quad = \boxed{\int (W_\varepsilon^2 \cdot P_0(dx))} = \int (\sqrt{\frac{P_\varepsilon}{P_0}} - 1)^2 P_0 \mu(dx) \\
 &\quad = \int \left(\frac{P_\varepsilon}{P_0} + 1 - 2 \sqrt{\frac{P_\varepsilon}{P_0}} \right) P_0 d\mu \\
 &\quad = \int (P_\varepsilon + P_0 - 2 \sqrt{P_\varepsilon P_0}) d\mu \\
 &\quad = \int (\sqrt{P_\varepsilon} - \sqrt{P_0})^2 d\mu \\
 &\quad + 2 \int [\sqrt{P_\varepsilon} - \sqrt{P_0}]^2 r(W_\varepsilon) d\mu
 \end{aligned}$$

$$\textcircled{2} \int (\sqrt{P_\varepsilon} - \sqrt{P_0}) \sqrt{P_0} d\mu$$

2st

$$\int (\sqrt{P_\varepsilon} - \sqrt{P_0})^2 d\mu \quad H\text{-distance}$$

$$+ 2 \int (\sqrt{P_\varepsilon} - \sqrt{P_0})^2 r(W_\varepsilon) d\mu \quad \begin{matrix} H\text{-distance} \\ \text{small order term} \end{matrix}$$

$$\begin{aligned} (1st) &= 2 \int (\sqrt{P_\varepsilon P_0} - P_0) d\mu \\ &= -2 \int (P_0 - \sqrt{P_0 P_\varepsilon}) d\mu \\ &= -2 \int \frac{P_\varepsilon + P_0 - 2\sqrt{P_\varepsilon P_0}}{2} d\mu \\ &= -\frac{2}{2} \int (\sqrt{P_\varepsilon} - \sqrt{P_0})^2 d\mu \\ &= -H^2(P_\varepsilon, P_0) \end{aligned}$$

$$(2nd) = H^2(P_\varepsilon, P_0)$$

$$(3rd) = o(\varepsilon^2)$$

↑ left to the students



$$P_0 [m_{\theta_0 + \varepsilon h} - m_{\theta_0}]$$

$$= -2 H^2(P_\varepsilon, P_0) + o(\varepsilon^2)$$

Wish to show $\frac{1}{2} \varepsilon^2 h^\top V_0 h$

It remains to show

$-2 H^2(P_\varepsilon, P_0)$ and $\frac{1}{2} \varepsilon^2 h^\top V_0 h$ (**) are $o(\varepsilon^2)$ close.

To show (*) holds:

① Reverse triangle inequality:

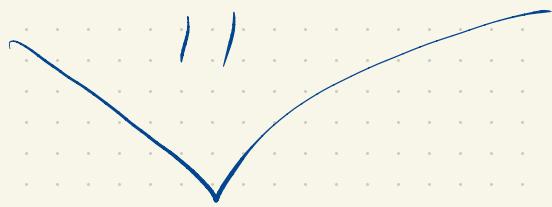
$$|||a|| - ||b||| \leq ||a - b||$$

$$(||a|| \leq ||b|| + ||a - b||)$$

$$\Rightarrow ||a|| - ||b|| \leq ||a - b||$$

② Introduce $L^2(\mu) = \{f : \frac{\|f\|_{L_2(\mu)}}{\left[\int f^2(x) \mu(dx) \right]^{1/2}} < \infty\}$

is a metric space



$$\begin{aligned} & \left| \| \sqrt{P_\Sigma} - \sqrt{P_0} \|_{L_2(\mu)} - \left\| \frac{1}{2} \varepsilon h^\top \dot{\ell}_{\theta_0} \sqrt{P_0} \right\|_{L_2(\mu)} \right| \\ & \leq \boxed{\left| \| \sqrt{P_\Sigma} - \sqrt{P_0} - \frac{1}{2} \varepsilon h^\top \dot{\ell}_{\theta_0} \sqrt{P_0} \|_{L_2(\mu)} \right|} \end{aligned}$$

QMD says $\square^2 = o(\varepsilon^2)$

which means

$$\| \sqrt{P_\Sigma} - \sqrt{P_0} \|_{L_2(\mu)} = \frac{1}{2} \varepsilon \| h^\top \dot{\ell}_{\theta_0} \sqrt{P_0} \|_{L_2(\mu)} + o(\varepsilon)$$

$$\boxed{\| \sqrt{P_\Sigma} - \sqrt{P_0} \|_{L_2(\mu)}^2} = \frac{1}{4} \varepsilon^2 \| h^\top \dot{\ell}_{\theta_0} \sqrt{P_0} \|_{L_2(\mu)}^2 + o(\varepsilon^2)$$

$H^2(P_\Sigma, P_0)$

$$-2H^2(P_\varepsilon, P_0) = -\frac{1}{2} \varepsilon^T \left[h^T \dot{\ell}_{\theta_0} (\bar{P}_0) \right]_{L_2(\mu)} + o(\varepsilon^2)$$

$$\begin{aligned} & h^T \left[\int \dot{\ell}_{\theta_0} \dot{\ell}_{\theta_0}^T \bar{P}_0 d\mu \right] h \\ & h^T \bar{P}_0 \dot{\ell}_{\theta_0} \dot{\ell}_{\theta_0}^T h \end{aligned}$$

$$P_0 [m_{\theta_0 + \varepsilon h} - m_{\theta_0}] = \frac{1}{2} \varepsilon^2 h^T V_0 h + o(\varepsilon^2)$$

$$\text{with } V_0 = P_0 [\dot{\ell}_{\theta_0} \dot{\ell}_{\theta_0}^T]$$

[Thm 7.6 rdV] Suppose

$\forall \theta \in$ open subset of \mathbb{R}^d ,
we have

$$P_\theta = \frac{dP_\theta}{d\mu}$$

Assume

(i) $\theta \mapsto \sqrt{P_\theta}$ is cont. differentiable $\forall x$,

(ii) The elements in the FIM

$$I(\theta) = \int \frac{\dot{P}_\theta(x) \dot{P}_\theta(x)^T}{[\dot{P}_\theta(x)]^2} P_\theta(dx)$$

are all well-defined and cont. w.r.t. θ

Then ^(a) $\theta \mapsto \overline{P\theta}$ is QMD.

(b) i_θ in the QMD is $\frac{\dot{P}_\theta}{P_\theta}$.