

## Outline

- 1) Maximum Likelihood Estimator
- 2) Asymptotic Distribution of MLE
- 3) Consistency of MLE

## Maximum Likelihood Estimation

For a generic dominated family  $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$  with densities  $p_\theta$ , a simple estimator for  $\theta$  is

$$\hat{\theta}_{MLE}(x) = \underset{\theta \in \Theta}{\operatorname{argmax}} p_\theta(x)$$

$$= \underset{\theta \in \Theta}{\operatorname{argmax}} l(\theta; x)$$

Remark 1:  $\operatorname{argmax}$  may not exist, be unique, or be computable

Remark 2: doesn't depend on parameterization or base measure, MLE for  $g(\theta)$  is  $g(\hat{\theta}_{MLE})$

$$\underline{Ex} \quad p_\gamma(x) = e^{\gamma' T(x) - A(\gamma)} h(x)$$

$$l(\gamma; x) = \gamma' T(x) - A(\gamma) + \log h(x)$$

$$\nabla l(\gamma; x) = T(x) - \mathbb{E}_\gamma T(x)$$

$$\Rightarrow \hat{\gamma}_{MLE} \text{ solves } T = \mathbb{E}_{\hat{\gamma}} T \quad \text{if such } \gamma \text{ exists}$$

Because  $\nabla^2 l(\gamma; x) = -\text{Var}_\gamma(T)$  is negative definite unless  $\gamma' T \stackrel{a.s.}{=} 0$  (in which case param. redundant)  
 $\Rightarrow$  at most 1 solution exists

Let  $\chi(\cdot)$  inverse of  $\mu(\gamma) = \nabla A(\gamma)$ ,  $\hat{\gamma} = \chi(T)$

$$\underline{\text{Ex}} \quad X_i \stackrel{\text{iid}}{\sim} e^{\gamma T(x) - A(\gamma)} h(x) \quad \gamma \in \Xi \subseteq \mathbb{R}$$

$$\hat{\gamma} = \psi(\bar{T}), \quad \bar{T} = \frac{1}{n} \sum T(x_i)$$

Assume  $\gamma \in \Xi^\circ$ .  $\dot{u}(\gamma) = \ddot{A}(\gamma) > 0 \quad \forall \gamma \in \Xi^\circ$

$$\text{so } \psi \text{ cts, } \dot{\psi}(u(\gamma)) = \frac{1}{\dot{u}(\gamma)} = \frac{1}{\ddot{A}(\gamma)}$$

Consistency:  $\bar{T} \xrightarrow{P} \mu(\gamma)$

Cts mapping:  $\psi(\bar{T}) \xrightarrow{P} \psi(\mu(\gamma)) = \gamma$

$$\text{Since } \sqrt{n}(\bar{T} - \mu(\gamma)) \Rightarrow N(0, \text{Var}_n(T(x_i))) \\ = N(0, \ddot{A}(\gamma)),$$

Delta method:

$$\sqrt{n}(\hat{\gamma} - \gamma) = \sqrt{n}(\psi(\bar{T}) - \psi(\mu(\gamma))) \\ \Rightarrow N(0, \frac{1}{\dot{A}(\gamma)^2} \cdot \ddot{A}(\gamma)) \\ = N(0, \frac{1}{\ddot{A}(\gamma)})$$

Recall  $J_1(\gamma) = \text{Var}_{\gamma}(T(x_i)) = \ddot{A}(\gamma)$   
 $= \text{Fisher info from 1 obs}$

$$\hat{\gamma} \approx N(\gamma, \frac{1}{n J_1(\gamma)})$$

Asymptotically unbiased, Gaussian, achieves CRLB

Ex  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Pois}(\theta)$ ,  $\gamma = \log \theta$

$$\hat{\gamma} = \log \bar{X}, \sqrt{n}(\bar{X} - \theta) \Rightarrow N(0, \theta)$$

$$\begin{aligned}\sqrt{n}(\hat{\gamma} - \gamma) &= \sqrt{n}(\log \bar{X} - \log \theta) \\ &\Rightarrow N(0, \theta \cdot \frac{1}{\theta^2}) \\ &= N(0, \theta^{-1})\end{aligned}$$

But If finite  $n$ ,  $\forall \theta > 0$ :

$$\begin{aligned}P_\theta(\hat{\gamma} = -\infty) &= P_\theta(X_i = 0)^n \\ &= e^{-\theta n} \rightarrow 0\end{aligned}$$

$$\Rightarrow E \hat{\gamma} = -\infty \quad \text{Var}(\hat{\gamma}) = \infty$$

[MLE can have embarrassing finite-sample performance despite being asy. optimal!]

Prop: If  $P(B_n) \rightarrow 0$ ,  $X_n \Rightarrow X$ ,  $Z_n$  arbitrary  
then  $X_n 1_{B_n^c} + Z_n 1_{B_n} \Rightarrow X$

Proof  $P(\|Z_n 1_{B_n}\| > \varepsilon) \leq P(B_n) \rightarrow 0$  so  $Z_n 1_{B_n} \xrightarrow{P} 0$

Also  $1_{B_n^c} \xrightarrow{P} 1$ , apply Slutsky  $\square$

[So zany behavior has no effect on cug. in dist]

## Asymptotic Efficiency

[The nice behavior of MLE we found in the exponential family case generalizes to a much broader class of models]

Setting  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} p_\theta(x) \quad \theta \in \mathbb{R}^d$

$p_\theta$  "smooth" in  $\theta$ , e.g. 2 cts integrable deriv.s  
(can be relaxed)

Let  $l_i(\theta; X_i) = \log p_\theta(X_i)$ ,  $l_n(\theta; X) = \sum_{i=1}^n l_i(\theta; X_i)$

$$J_1(\theta) = \text{Var}_\theta(\nabla l_1(\theta; X_i)) = -\mathbb{E}_\theta[\nabla^2 l_1(\theta; X_i)]$$

$$J_n(\theta) = \text{Var}_\theta(\nabla l_n(\theta; X)) = n J_1(\theta)$$

We say an estimator  $\hat{\theta}_n$  is asymptotically efficient if  $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{P_\theta} N(0, J_1(\theta)^{-1})$

Delta method for estimand  $g(\theta)$ :

$$\sqrt{n}(g(\hat{\theta}_n) - g(\theta)) \xrightarrow{P_\theta} N(0, \nabla g(\theta)' J_1(\theta)^{-1} \nabla g(\theta))$$

also achieves CRLB if  $\hat{\theta}_n$  does;  $g$  diff.

## Asymptotic Dist. of MLE

Under mild conditions,  $\hat{\theta}_{MLE}$  is asy. Gaussian, efficient

We will be interested in  $\ell(\theta; X)$  as a function of  $\theta$

Note "true" value as  $\theta_0$  ( $X \sim P_{\theta_0}$ )

Derivatives of  $\ell_n$  at  $\theta_0$ :

$$\nabla \ell_1(\theta_0; X_i) \stackrel{iid}{\sim} (0, J_1(\theta_0))$$

$$\frac{1}{\sqrt{n}} \nabla \ell_n(\theta_0; X) = \sqrt{n} \cdot \frac{1}{n} \sum \nabla \ell_1(\theta_0; X_i) \xrightarrow{P_{\theta_0}} N(0, J_1(\theta_0))$$

$$\frac{1}{n} \nabla^2 \ell_n(\theta_0; X) \xrightarrow{P_{\theta_0}} \mathbb{E}_{\theta_0} \nabla^2 \ell_1(\theta_0; X_i) = -J_1(\theta_0)$$

Informal Proof:

$$0 = \nabla \ell_n(\hat{\theta}_n; X) \approx \nabla \ell_n(\theta_0) + \nabla^2 \ell_n(\theta_0)(\hat{\theta}_n - \theta_0)$$

$$\begin{aligned} \sqrt{n}(\hat{\theta}_n - \theta_0) &\approx - \underbrace{\left( \frac{1}{n} \nabla^2 \ell_n(\theta_0) \right)^{-1}}_{\xrightarrow{P} J(\theta_0)^{-1}} \underbrace{\frac{1}{\sqrt{n}} \nabla \ell_n(\theta_0)} \\ &\Rightarrow N_d(0, J(\theta_0)^{-1}) \end{aligned}$$

More rigorous proof later, but note  
we need consistency of  $\hat{\theta}_n$  first to even  
justify Taylor expansion

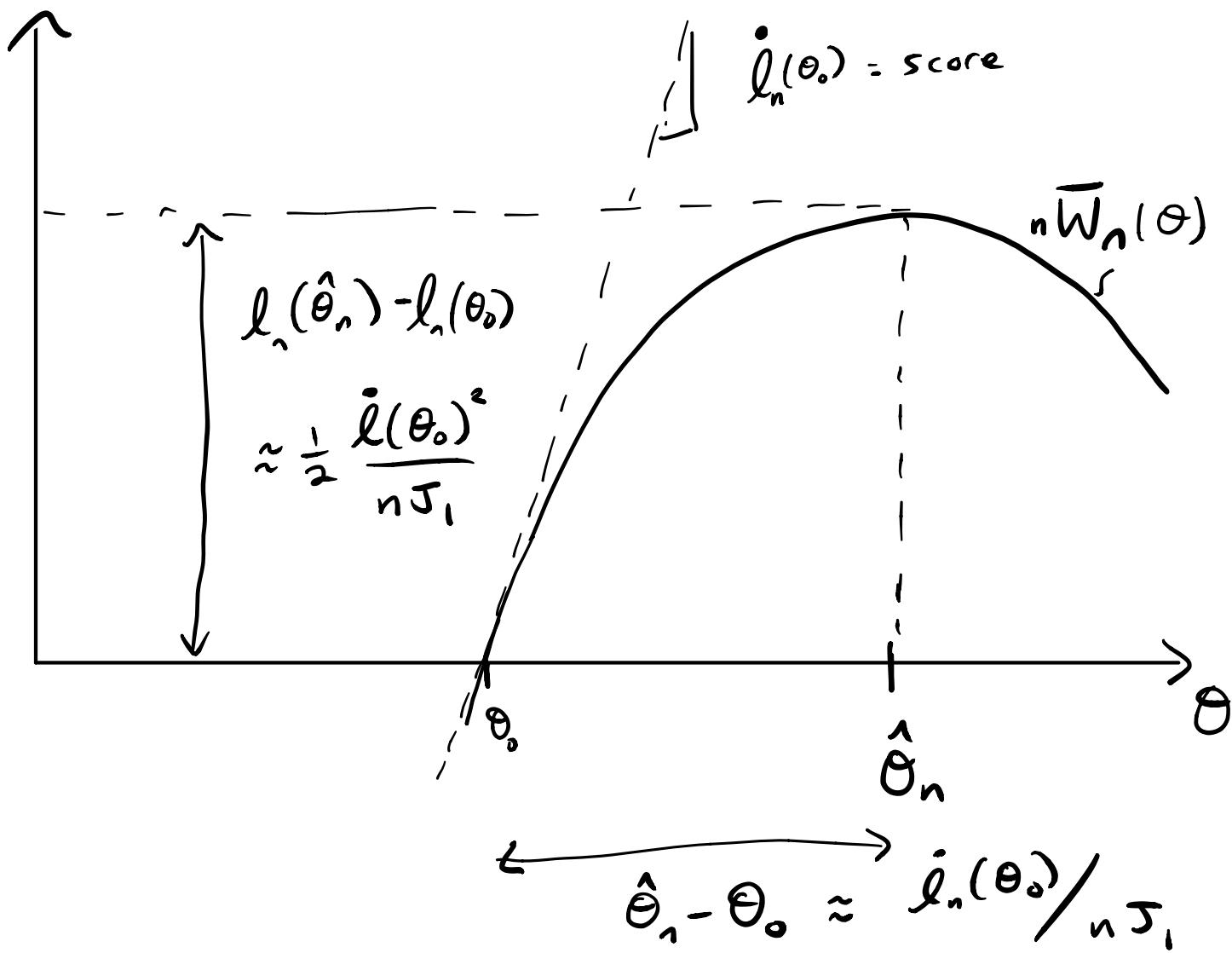
# Asymptotic Picture ( $d=1$ )

Recall  $(l_n(\theta) - l_n(\theta_0))_{\theta \in \mathbb{U}}$  is minimal suff.

Quadratic approximation near  $\theta_0$ :

$$l_n(\theta) - l_n(\theta_0) \approx \underbrace{\dot{l}_n(\theta_0)}_{\approx N(0, nJ_1(\theta_0))} (\theta - \theta_0) + \frac{1}{2} \underbrace{\ddot{l}_n(\theta_0)}_{\approx -nJ_1(\theta_0)} (\theta - \theta_0)^2$$

Gaussian linear term                              Deterministic curvature



## Consistency of MLE

$$X_1, \dots, X_n \stackrel{iid}{\sim} P_{\theta_0}, \quad \hat{\theta}_n \in \operatorname{argmax}_{\theta \in \Theta} l_n(\theta; X)$$

[Assuming  $P(\operatorname{argmax} \neq \emptyset) \rightarrow 1$  as  $n \rightarrow \infty$ ]

Question: When does  $\hat{\theta}_n \xrightarrow{P} \theta_0$  ?

Assume model identifiable ( $P_\theta \neq P_{\theta_0}$  for  $\theta \neq \theta_0$ )

Recall KL Divergence:

$$\begin{aligned} D_{KL}(\theta_0 \parallel \theta) &= E_{\theta_0} \log \frac{P_{\theta_0}(X_i)}{P_\theta(X_i)} \\ -D_{KL}(\theta_0 \parallel \theta) &\leq \log E_{\theta_0} \frac{P_\theta(X_i)}{P_{\theta_0}(X_i)} \quad \leftarrow \text{(note switch)} \\ &= \log \int_{x: P_{\theta_0}(x) > 0} \frac{P_\theta(x)}{P_{\theta_0}(x)} P_{\theta_0}(x) d\mu(x) \\ &\leq \log 1 = 0 \end{aligned}$$

(Jensen)

strict inequality unless  $\frac{P_\theta}{P_{\theta_0}} \text{ const.}$  (i.e., unless  $P_\theta = P_{\theta_0}$ )

Let  $W_i(\theta) = l_i(\theta; X_i) - l(\theta_0; X_i)$ ,  $\bar{W}_n = \frac{1}{n} \sum W_i$

Note  $\hat{\theta}_n \in \operatorname{argmax}_{\theta \in \Theta} \bar{W}_n(\theta)$  too

$$\begin{aligned}
 \bar{W}_n(\theta) &\xrightarrow{P} \mathbb{E}_{\theta_0} W_i(\theta) \\
 &= -D_{KL}(\theta_0 \parallel \theta) \\
 &\leq 0, \text{ equality iff } \theta = \theta_0
 \end{aligned}$$

But not enough: need uniform convergence in  $\theta$

Def For compact  $K$  let  $C(K) = \{f: K \rightarrow \mathbb{R}, \text{cts}\}$

For  $f \in C(K)$  let  $\|f\|_\infty = \sup_{t \in K} |f(t)|$

$f_n \rightarrow f$  in this norm if  $\|f_n - f\|_\infty \rightarrow 0$

Thm (LLN for random functions)

Assume  $K$  compact,  $W_1, W_2, \dots \in C(K)$  iid.

$$\mathbb{E} \|W_i\|_\infty < \infty, \quad \mu(t) = \mathbb{E} W_i(t)$$

Then  $\mu(t) \in C(K)$

and  $P\left(\|\frac{1}{n} \sum W_i - \mu\|_\infty > \varepsilon\right) \rightarrow 0$

(i.e.,  $\bar{W}_n \xrightarrow{P} \mu$  in  $\|\cdot\|_\infty$ , or  $\|\bar{W}_n - \mu\|_\infty \xrightarrow{P} 0$ )

## Theorem (Keener 9.4):

Let  $G_1, G_2, \dots$  random functions in  $C(K)$ ,  $K$  cpt.

$\|G_n - g\|_\infty \xrightarrow{P} 0$ , some fixed  $g \in C(K)$ . Then

① If  $t_n \xrightarrow{P} t^* \in K$  ( $t^*$  fixed) then  $G_n(t_n) \xrightarrow{P} g(t^*)$

② If  $g$  maximized at unique value  $t^*$ ,

and  $G_n(t_n) = \max G_n(t)$  then  $t_n \xrightarrow{P} t^*$   
 $G_n(t_n) \geq \max G_n - \alpha_n, \alpha_n \rightarrow 0$  (mod.s of proof in purple)

③ If  $K \subseteq \mathbb{R}$ ,  $g(t) = 0$  has unique sol.  $t^*$ ,

and  $t_n$  solve  $G_n(t_n) = 0$  then  $t_n \xrightarrow{P} t^*$

$$|G_n(t_n)| \leq \alpha_n, \alpha_n \rightarrow 0$$

## Proof

$$\begin{aligned} ① |G_n(t_n) - g(t^*)| &\leq |G_n(t_n) - g(t_n)| + |g(t_n) - g(t^*)| \\ &\leq \|G_n - g\|_\infty + |g(t_n) - g(t^*)| \\ &\xrightarrow{P} 0 \quad \xrightarrow{P} 0 \\ &\text{(by assumptions)} \quad \text{(by cts mapping)} \end{aligned}$$

② Fix  $\varepsilon > 0$ , let  $B_\varepsilon(t^*) = \{t : \|t - t^*\| < \varepsilon\}$

Let  $K_\varepsilon = K \setminus B_\varepsilon(t^*) = K \cap B_\varepsilon^c(t^*)$  (compact)

$$\delta = g(t^*) - \max_{t \in K_\varepsilon} g(t) > 0$$

If  $t_n \in K_\varepsilon$  then  $G_n(t_n) \leq g(t^*) - \delta + \|G_n - g\|_\infty$

and  $G_n(t_n) \geq G_n(t^*) \geq g(t^*) - \|G_n - g\|_\infty$

then  $2\|G_n - g\|_\infty \geq \delta$

$P(\|t_n - t^*\| \geq \varepsilon) \leq P(\|G_n - g\|_\infty \geq \frac{\delta}{2}) \rightarrow 0$

③ Analogous to ②

Theorem (Consistency of MLE for compact  $\Theta$ )

$X_1, \dots, X_n \stackrel{iid}{\sim} P_{\theta_0}$ ,  $\mathcal{P}$  has cts densities  $P_\theta$ ,  $\Theta \in \mathbb{H}$

Assume

- $\Theta$  compact
- $\mathbb{E}_{\Theta_0} \|W\|_\infty = \mathbb{E}_{\Theta_0} \|\ell_1(\theta; x_i) - \ell_1(\theta_0; x_i)\|_\infty < \infty$
- Model identifiable

Then  $\hat{\theta}_n \xrightarrow{P} \theta_0$  if  $\hat{\theta}_n \in \arg\max \ell_n(\theta; X)$

Proof  $W_i \in C(\Theta)$  iid, mean  $\mu(\theta) = -D_{KL}(\theta_0 \parallel \theta)$

$\mu(\theta_0) = 0$ ,  $\mu(\theta) < 0 \quad \forall \theta \neq \theta_0$  ( $\theta_0 = \operatorname{arg\min} \mu$ )

By definition,  $\hat{\theta}_n$  maximizes  $\bar{W}_n$ ,

$\|\bar{W}_n - \mu\|_\infty \xrightarrow{P} 0$ , apply 9.4, ②

We usually care about non-compact parameter spaces, need some extra assumption to get us there.

Thm ( $\approx$  Keener 9.11, but stronger conditions)

$X_1, \dots, X_n \stackrel{iid}{\sim} p_{\theta_0}$ ,  $\mathcal{P}$  has cts densities  $p_{\theta}$ ,  $\theta \in \Theta = \mathbb{R}^d$

Assume . Model identifiable

- For all compact  $K \subseteq \mathbb{R}^d$ ,  $\mathbb{E} \left[ \sup_{\theta \in K} |W_i(\theta)| \right] < \infty$
- $\exists r > 0$  s.t.  $\mathbb{E} \left[ \sup_{\|\theta - \theta_0\| \geq r} W_i(\theta) \right] < 0$

Then  $\hat{\theta}_n \xrightarrow{P} \theta_0$  if  $\hat{\theta}_n \in \operatorname{argmax} l_n(\theta; X)$

Proof Let  $A = \{\theta : \|\theta - \theta_0\| > r\}$ ,  $\alpha = \mathbb{E} \sup_{\theta \in A} W_i(\theta) < 0$

$$\sup_{\theta \in A} \bar{W}_n(\theta) \leq \frac{1}{n} \sum_{i=1}^n \sup_{\theta \in A} W_i(\theta) \rightarrow \alpha < 0$$

Hence,  $P(\hat{\theta}_n \in A) \leq P(\underbrace{\bar{W}_n(\theta_0)}_{\xrightarrow{P} 0} < \underbrace{\sup_{\theta \in A} \bar{W}_n(\theta)}_{\xrightarrow{P} \alpha}) \rightarrow 0$

Let  $\hat{\theta}_n^A = \hat{\theta}_n \mathbf{1}\{\hat{\theta}_n \in A^c\} + \theta_0 \mathbf{1}\{\hat{\theta}_n \in A\}$

$\xrightarrow{P} \theta_0$  by previous thm., since  $A^c$  is compact

Hence  $\hat{\theta}_n \xrightarrow{P} \theta_0$  by our Proposition.

# Asymptotic Dist. of MLE

## Theorem

$X_1, \dots, X_n \stackrel{iid}{\sim} P_{\theta_0}$  for  $\theta_0 \in \Theta^0 \subseteq \mathbb{R}^d$

Assume •  $\hat{\theta}_n \in \underset{\theta \in \Theta}{\operatorname{argmax}} l_n(\theta; x)$ ,  $\hat{\theta}_n \xrightarrow{P} \theta_0$

• In a neighborhood  $B_\varepsilon(\theta_0) = \{\theta : \|\theta - \theta_0\| < \varepsilon\} \subseteq \Theta$ :

(i)  $l_1(\theta; x)$  has 2 cts deriv.s on  $B_\varepsilon(\theta_0)$ ,  $\forall x$

(ii)  $\mathbb{E}_{\theta_0} \left[ \sup_{\theta \in B_\varepsilon} \|\nabla^2 l_1(\theta; X_i)\| \right] < \infty$

any norm on  $\mathbb{R}^{d \times d}$ ,  
e.g. Frobenius

• Fisher info :

$$\mathbb{E}_{\theta_0} \nabla l_1(\theta_0; X_i) = 0$$

$$\text{Var}_{\theta_0} \nabla l_1(\theta_0; X) = -\mathbb{E}_{\theta_0} \nabla^2 l_1(\theta_0; X) > 0$$

(enough to have 3<sup>rd</sup> deriv. of  $l_1$  bdd in  $B_\varepsilon(\theta_0)$ )

Then  $\sqrt{n}(\hat{\theta}_n - \theta_0) \Rightarrow N_d(0, J_1(\theta_0)^{-1})$

Proof Let  $A_n = \{ \|\hat{\theta}_n - \theta_0\| \geq \varepsilon \}$ ,

$$P_{\theta_0}(A_n) \rightarrow 0 \quad \text{by assumption}$$

On  $A_n^c$ ,  $\hat{\theta}_n \in B_\varepsilon(\theta_0)$  and we have

$$0 = \nabla l_n(\hat{\theta}_n; x)$$

$$= \nabla l_n(\theta_0; x) + \nabla^2 l_n(\tilde{\theta}_n; x)(\hat{\theta}_n - \theta_0),$$

for some  $\tilde{\theta}_n$  between  $\theta_0$  and  $\hat{\theta}_n$  (MVT)

$$\begin{aligned} \sqrt{n}(\hat{\theta}_n - \theta_0) &= \underbrace{\left( -\frac{1}{n} \nabla^2 l_n(\tilde{\theta}_n) \right)^{-1}}_{\xrightarrow{\text{P.s.}} J_1(\theta_0)^{-1}} \underbrace{\frac{1}{\sqrt{n}} \nabla l_n(\theta_0)}_{\xrightarrow{\text{By 9.14 ① + cts mapping}} N(0, J_1(\theta_0))} \\ &\Rightarrow N_d(0, J_1(\theta_0)^{-1}) \end{aligned}$$

$$\Rightarrow N_d(0, J_1(\theta_0)^{-1})$$

Behavior on  $A_n$  irrelevant to asymptotic limit  $\square$

[More rigorous: define "fixed" estimator

$$\bar{\theta}_n = \begin{cases} \hat{\theta}_n & \text{on } A_n^c \\ \theta_0 + J_1(\theta_0)^{-1} \frac{1}{n} \nabla l_n(\theta_0) & \text{on } A_n \end{cases}$$

$$\begin{aligned} \sqrt{n}(\bar{\theta}_n - \theta_0) &= \left( J_1 \mathbf{1}_{A_n} + \frac{1}{n} \nabla^2 l(\tilde{\theta}_n) \mathbf{1}_{A_n^c} \right)^{-1} \frac{1}{\sqrt{n}} \nabla l_n(\theta_0) \\ &\Rightarrow N_d(0, J_1^{-1}) \end{aligned}$$

$$\hat{\theta}_n = \bar{\theta}_n \mathbf{1}_{A_n^c} + \hat{\theta}_n \mathbf{1}_{A_n}$$

]