

11/4/2021

# 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

$$\begin{aligned}\hat{\theta}_{MLE}(x) &= \operatorname{argmax}_{\theta \in \Theta} p_\theta(x) \\ &= \operatorname{argmax}_{\theta \in \Theta} \ell(\theta; x)\end{aligned}$$

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  $g(\hat{\theta}_{MLE})$

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

$$\ell(\eta; x) = \eta' T(x) - A(\eta) + \log h(x)$$

$$\nabla \ell(\eta; x) = T(x) - \mathbb{E}_\eta T(x)$$

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

Because  $\nabla^2 \ell(\eta; x) = -\operatorname{Var}_\eta(T)$  is negative definite unless  $\eta' T \stackrel{\text{a.s.}}{=} 0$  (in which case param. redundant)

$\Rightarrow$  at most 1 solution exists

$$\text{Let } \mu = \eta(\eta) = \nabla A(\eta), \quad \hat{\eta} = \eta^{-1}(T)$$

(HW 4.2)

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

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

$$\text{Assume } \eta \in \Xi^o. \quad \dot{\psi}(\eta) = \ddot{A}(\eta) > 0 \quad \forall \eta \in \Xi^o$$

$$\text{so } \psi^{-1} \text{ cts, } (\dot{\psi}^{-1})(\mu) = \frac{1}{\dot{\psi}(\psi(\mu))} = \frac{1}{\ddot{A}(\eta)}$$

$$\text{Consistency: } \bar{T} \xrightarrow{P_n} \mu$$

$$\text{Cts mapping: } \psi^{-1}(\bar{T}) \xrightarrow{P_n} \psi^{-1}(\mu) = \eta$$

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

Delta method:

$$\text{(Recall } J_1(\mu) = \text{Var}(T)^{-1} \\ = \ddot{A}(\eta)^{-1} \text{)}$$

$$\sqrt{n}(\hat{\eta} - \eta) = \sqrt{n}(\psi^{-1}(\bar{T}) - \eta)$$

$$\Rightarrow N(0, \frac{1}{\ddot{A}(\eta)^2} \cdot \ddot{A}(\eta))$$

$$= N(0, \frac{1}{\ddot{A}(\eta)})$$

$$\text{Recall } J_1(\eta) = \text{Var}_{\eta}(T(x_i)) = \ddot{A}(\eta)$$

= Fisher info from 1 obs

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

Asymptotically unbiased, Gaussian, achieves CRLB

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

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

$$\sqrt{n}(\hat{\eta} - \eta) = \sqrt{n}(\log \bar{X} - \log \theta)$$

$$\Rightarrow N(0, \theta \cdot \frac{1}{\theta^2})$$

$$= N(0, \theta^{-1})$$

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

$$P_{\theta}(\hat{\eta} = -\infty) = P_{\theta}(X_1 = 0)^n$$

$$= e^{-\theta n} > 0$$

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

[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

$$\text{then } X_n 1_{B_n^c} + Z_n 1_{B_n} \Rightarrow X$$

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  $\boxtimes$

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

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

Let  $l_1(\theta; X_i) = \log p_\theta(X_i)$ ,  $l_n(\theta; X) = \sum_{i=1}^n l_1(\theta; X_i)$

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

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

We say an estimator  $\hat{\theta}_n$  is asymptotically efficient

$$\text{if } \sqrt{n}(\hat{\theta}_n - \theta) \stackrel{P_\theta}{\Rightarrow} \mathcal{N}(0, J_1(\theta)^{-1})$$

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

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

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  $l_n(\theta; X)$  as a function of  $\theta$

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

Derivatives of  $l_n$  at  $\theta_0$ :

$$\nabla l_1(\theta_0; X_i) \stackrel{iid}{\sim} (0, J_1(\theta_0))$$

$$\frac{1}{\sqrt{n}} \nabla l_n(\theta_0; X) = \sqrt{n} \cdot \frac{1}{n} \sum \nabla l_1(\theta_0; X_i) \xrightarrow{P_{\theta_0}} N(0, J_1(\theta_0))$$

$$\frac{1}{n} \nabla^2 l_n(\theta_0; X) \xrightarrow{P_{\theta_0}} E_{\theta_0} \nabla^2 l_1(\theta_0; X_i) = -J_1(\theta_0)$$

Informal Proof:

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

↙ between  $\theta_0, \tilde{\theta}_n$

$$\sqrt{n} (\hat{\theta}_n - \theta_0) = - \underbrace{\left( \frac{1}{n} \nabla^2 l_n(\tilde{\theta}_n) \right)^{-1}}_{\text{(want)}} \underbrace{\frac{1}{\sqrt{n}} \nabla l_n(\theta_0)}_{\Rightarrow N_d(0, J(\theta_0))}$$

$$\begin{aligned} &\text{(want)} \xrightarrow{P} J(\theta_0)^{-1} \Rightarrow N_d(0, J(\theta_0)) \\ &\Rightarrow N_d(0, J(\theta_0)^{-1}) \end{aligned}$$

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

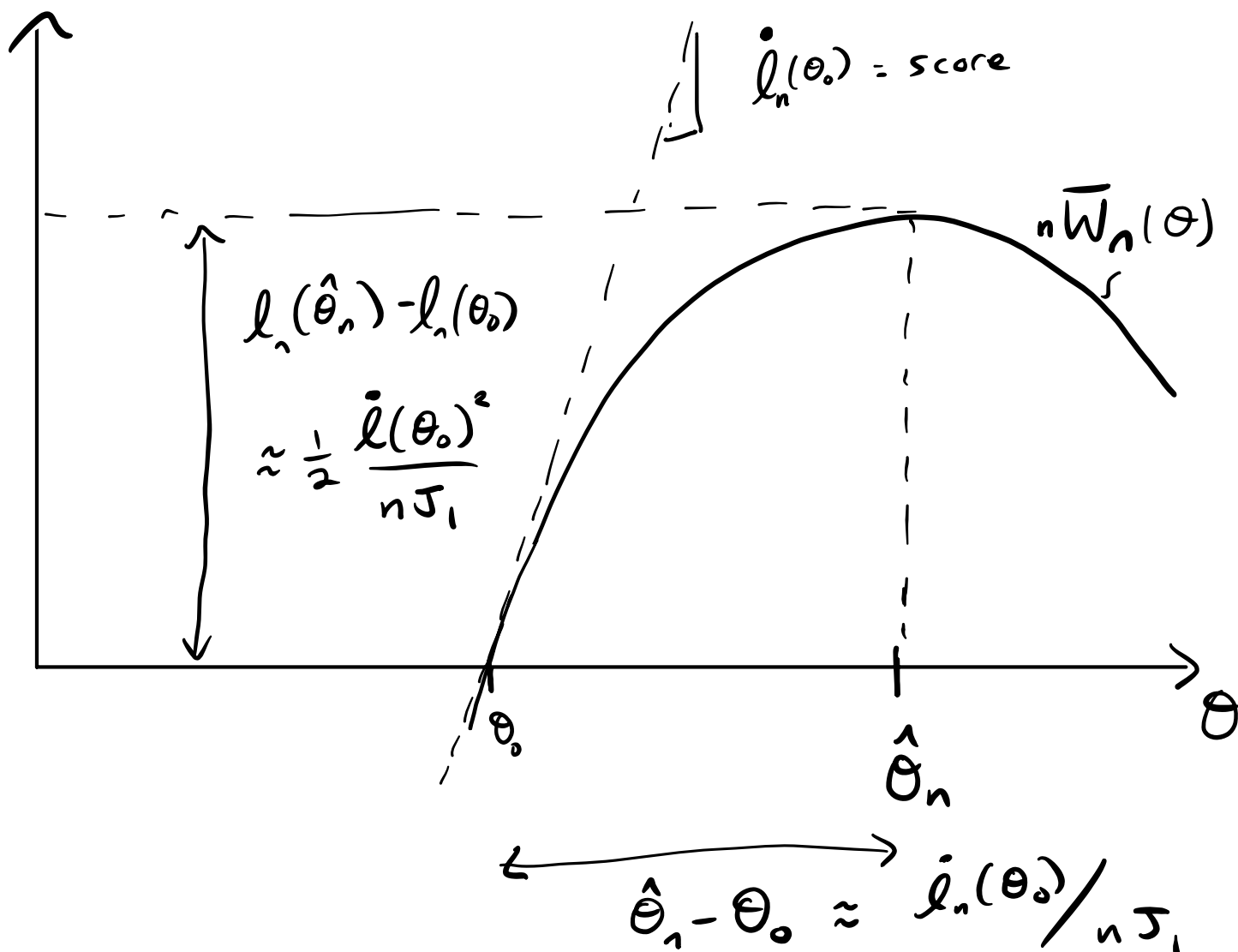
# Asymptotic Picture (d=1)

Recall  $(l_n(\theta) - l_n(\theta_0))_{\theta \in \Theta}$  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

$$X_1, \dots, X_n \stackrel{iid}{\sim} P_{\theta_0}, \quad \hat{\theta}_n \in \arg \max_{\theta \in \Theta} \ell_n(\theta; X)$$

[ Will be ok if  $\hat{\theta}_n$  comes close to maximizing  $\ell_n$  ]

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:

$$D_{KL}(\theta_0 \parallel \theta) = \mathbb{E}_{\theta_0} \log \frac{P_{\theta_0}(X_i)}{P_{\theta}(X_i)}$$

$$-D_{KL}(\theta_0 \parallel \theta) \leq \log \mathbb{E}_{\theta_0} \frac{P_{\theta}(X_i)}{P_{\theta_0}(X_i)} \quad \leftarrow \text{(note switch)}$$

$$= \log \int \frac{P_{\theta}(x)}{P_{\theta_0}(x)} P_{\theta_0}(x) d\mu(x)$$

$$\leq \log 1 = 0$$

(Jensen)

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

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

Note  $\hat{\theta}_n \in \arg \max_{\theta \in \Theta} \bar{W}_n(\theta)$  too



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

But not enough:

- MLE  $\hat{\theta}_n$  depends on entire function  $\bar{W}_n(\cdot)$
- 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.

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

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

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

(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. 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

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

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  $G_n(t_n) \leq \underbrace{g(t^*) - \delta}_{> \max_{K_\varepsilon} g(t)} + \|G_n - g\|_\infty$

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

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

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

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  $p_\theta$ ,  $\theta \in \Theta$

Assume  $\cdot \Theta$  compact

$$\cdot \mathbb{E}_{\theta_0} \|W_i\|_\infty = \mathbb{E}_{\theta_0} \|\ell_1(\theta; X_i) - \ell_1(\theta_0; X_i)\|_\infty < \infty$$

$\cdot$  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  $\mu(\theta) = -D_{KL}(\theta_0 \| \theta)$

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

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

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

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  $\cdot$  Model identifiable

$\cdot$  For all compact  $K \subseteq \mathbb{R}^d$ ,  $\mathbb{E} \left[ \sup_{\theta \in K} |W_i(\theta)| \right] < \infty$

$\cdot \exists r > 0$  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  $\hat{\theta}_n \in \operatorname{argmax} \ell_n(\theta; X)$

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

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

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

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

$\xrightarrow{P} \theta_0$  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  $\theta_0 \in \Theta^o \subseteq \mathbb{R}^d$

Assume  $\cdot \hat{\theta}_n \in \operatorname{argmax}_{\theta \in \Theta} l_n(\theta; X)$ ,  $\hat{\theta}_n \xrightarrow{P} \theta_0$

$\cdot$  In a neighborhood  $\bar{B}_\varepsilon(\theta_0) = \{\theta : \|\theta - \theta_0\| \leq \varepsilon\} \subseteq \Theta^o$ :

(i)  $l_1(\theta; X)$  has 2 cts deriv.s on  $\bar{B}_\varepsilon(\theta_0)$ ,  $\forall X$

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

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

$\cdot$  Fisher info:

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

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

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

(positive def.)  
 $\downarrow$   
 $\checkmark \checkmark \checkmark \checkmark$

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

Proof  $\sup_{\theta \in \bar{B}_\varepsilon} \|\frac{1}{n} \nabla^2 \ell_n(\theta) - J_1(\theta)\| \xrightarrow{P} 0$

since  $\bar{B}_\varepsilon$  compact,  
 $\nabla^2 \ell_1$  cts

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

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

$O_n A_n^c$ ,  $\hat{\theta}_n \in \bar{B}_\varepsilon(\theta_0)$  and we have

$$0 = \nabla \ell_n(\hat{\theta}_n; X)$$

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

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

$$\sqrt{n} (\hat{\theta}_n - \theta_0) = \underbrace{\left( -\frac{1}{n} \nabla^2 \ell_n(\tilde{\theta}_n) \right)^{-1}}_{\xrightarrow{P_{\theta_0}} J_1(\theta_0)^{-1}} \underbrace{\frac{1}{\sqrt{n}} \nabla \ell_n(\theta_0)}_{\Rightarrow N(0, J_1(\theta_0))}$$

By 9.4 ① + cts mapping

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

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