

# Outline

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

# 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  $\theta_0$ : ( $\theta_0 \in \Theta^\circ$ )

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

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

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

Proof sketch:

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

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

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

$$\xrightarrow{P} J(\theta_0)^{-1} \Rightarrow N_d(0, J(\theta_0))$$

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

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

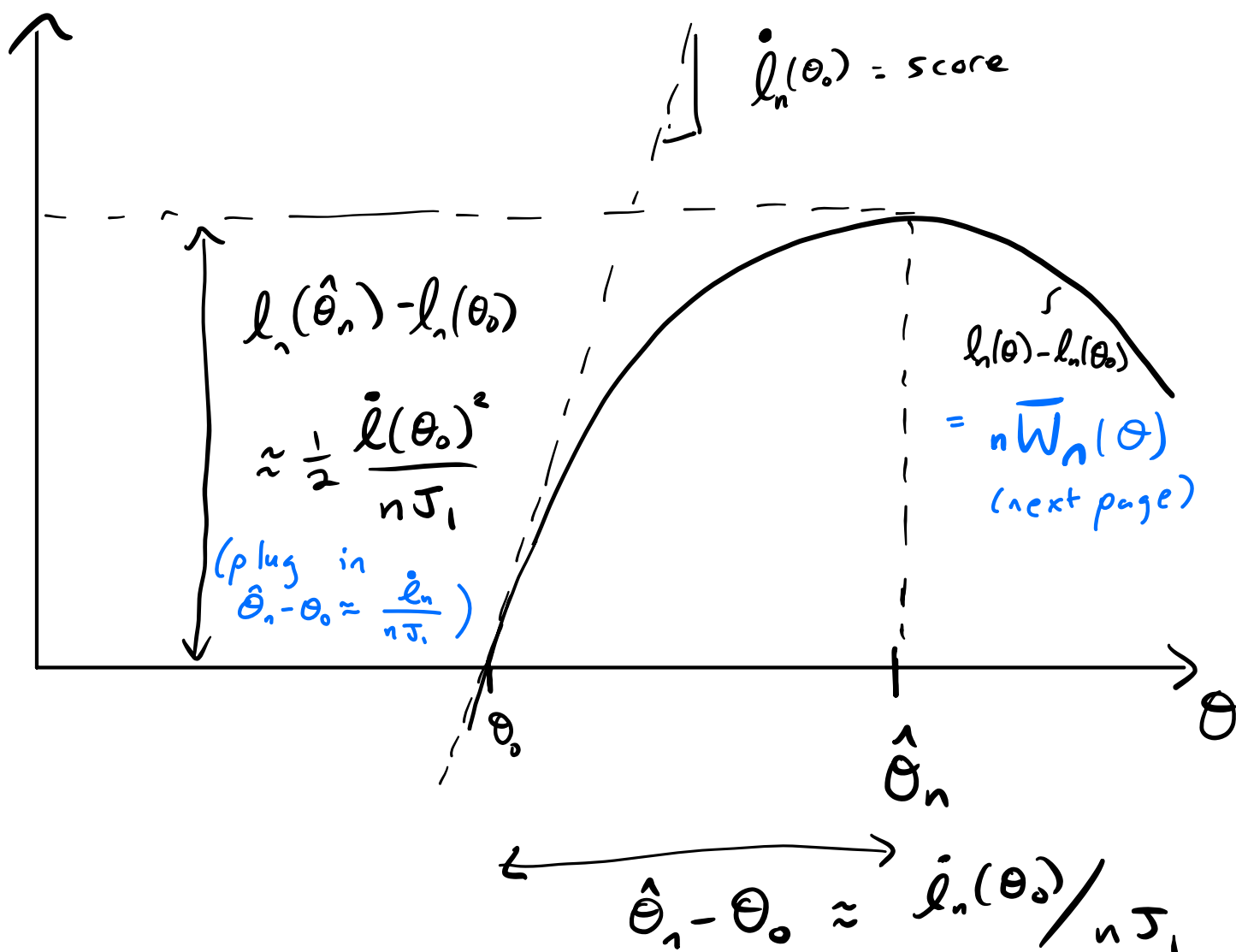
# Asymptotic Picture (d=1)

Recall  $(\ell_n(\theta) - \ell_n(\theta_0))_{\theta \in \Theta}$  is minimal suff.

Quadratic approximation near  $\theta_0$ :

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

Gaussian linear term
Deterministic curvature



# Consistency of MLE

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

[ 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  $\hat{\theta}_n \in \underset{\theta \in \Theta}{\operatorname{argmax}} \bar{W}_n(\theta)$  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  $\|f\|_{\infty} = \sup_{t \in K} |f(t)|$

$f_n \xrightarrow{P} f$  in this norm if  $\|f_n - f\|_{\infty} \xrightarrow{P} 0$

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 (Consistency of MLE for compact $\Theta$ )

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

- Assume
- $\log p_{\theta}(x)$  cts in  $\Theta$ , all  $x \in \mathcal{X}$
  - $\Theta$  compact
  - $\mathbb{E}_{\theta_0} \left[ \sup_{\theta \in \Theta} |W_i(\theta)| \right] < \infty = \mathbb{E}_{\theta_0} \sup_{\theta} |l(\theta; X_i) - l(\theta_0; X_i)|$
  - Model identifiable

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

Proof  $W_i \in C(\Theta)$  iid, mean  $\mu(\theta) = -D_{KL}(\theta_0 \| \theta)$   
 $\mu(\theta_0) = 0$ ,  $\mu(\theta) < 0 \quad \forall \theta \neq \theta_0$  ( $\theta_0 = \arg \min \mu$ )

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

$$\mathcal{J}_n = \|\bar{w}_n - \mu\|_{\infty} \xrightarrow{P} 0.$$

Fix  $\varepsilon > 0$ , want to show  $\mathbb{P}_{\theta_0}(\|\theta - \theta_0\| \geq \varepsilon) \rightarrow 0$

Let  $\tilde{\Theta}_{\varepsilon} = \Theta \setminus B_{\varepsilon}(\theta_0) = \{\theta \in \Theta : \|\theta - \theta_0\| \geq \varepsilon\}$  (compact)

$$\text{Let } \mu_{\varepsilon}^* = \max_{\theta \in \tilde{\Theta}_{\varepsilon}} \mu(\theta) < 0 = \mu(\theta_0)$$

$$w_{\varepsilon}^* = \max_{\theta \in \tilde{\Theta}_{\varepsilon}} \bar{w}_n(\theta)$$

$$\begin{aligned} \mathbb{P}_{\theta_0}(\|\hat{\theta}_n - \theta_0\| \geq \varepsilon) &\leq \mathbb{P}_{\theta_0}(\underbrace{w_\varepsilon^*}_{\leq \mu_\varepsilon^* + \delta_n} \geq \underbrace{\bar{w}_n(\theta_0)}_{\geq -\delta_n}) \\ &\leq \mathbb{P}_{\theta_0}(2\delta_n \geq \underbrace{-\mu_\varepsilon^*}_{> 0}) \rightarrow 0 \quad \square \end{aligned}$$

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

Corollary Same assumptions except now  $\Theta = \mathbb{R}^d$ , <sup>(non-compact)</sup>  
but there's some  $R < \infty$  large enough so

$$\mathbb{P}_{\theta_0}(\|\hat{\theta}_n - \theta_0\| > R) \rightarrow 0$$

Then  $\hat{\theta}_n \xrightarrow{P} \theta_0$

Proof Let  $\tilde{\Theta} = \{\theta : \|\theta - \theta_0\| \leq R\}$ ,  $\tilde{\theta}_n = \arg \max_{\theta \in \tilde{\Theta}} p_\theta(x)$

Then  $\tilde{\theta}_n \rightarrow \theta_0$  by assumption

$$\mathbb{P}_{\theta_0}(\hat{\theta}_n \neq \tilde{\theta}_n) = \mathbb{P}_{\theta_0}(\hat{\theta}_n \notin \tilde{\Theta}) \rightarrow 0$$

$$\text{so } \hat{\theta}_n - \tilde{\theta}_n \xrightarrow{P} 0 \Rightarrow \hat{\theta}_n = \tilde{\theta}_n + (\hat{\theta}_n - \tilde{\theta}_n) \xrightarrow{P} \theta_0 \quad \square$$

So the only thing we actually need to worry about is if  $\hat{\theta}_n$  is extremely far away from  $\theta_0$  with non-negligible Prob.

## Theorem (Asymptotic distribution of MLE)

$X_1, \dots, X_n \stackrel{iid}{\sim} p_{\theta_0}$     $\mathcal{P}$  has densities  $p_{\theta}$ ,  $\theta \in \Theta$

Assume

- $\mathcal{P}$  identifiable

- $\Theta$  compact

- $\mathbb{E}_{\theta_0} \left[ \sup_{\theta \in \Theta} |W_i(\theta)| \right] < \infty$

- $l(\theta; x) = \log p_{\theta}(x)$  has two cts derivatives in  $\Theta$

- $\mathbb{E}_{\theta_0} \sup_{\theta \in \Theta} \|\nabla^2 l_i(\theta; x_i)\| < \infty$

- $J_1(\theta_0) = \mathbb{E}_{\theta_0} \nabla^2 l_i(\theta_0; x_i) \succ 0$  ↙ positive definite

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

Proof From before, we had

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = \left( -\frac{1}{n} \nabla^2 l_n(\tilde{\theta}_n) \right)^{-1} \nabla l_n(\theta_0)$$

for  $\tilde{\theta}_n$  between  $\theta_0$  and  $\hat{\theta}_n$

Previous result shows  $\hat{\theta}_n \xrightarrow{P} \theta_0$ , so  $\tilde{\theta}_n \xrightarrow{P} \theta_0$  also

Define  $V_i(\theta) = -\nabla^2 l_i(\theta; x_i) \in C(\Theta)$ ,  $\mathbb{E}_{\theta_0} \|V_i\|_{\infty} < \infty$   
by assumption

Then  $v(\theta) = \mathbb{E}_{\theta_0} V_i(\theta) \in C(\Theta)$ ,  $v(\theta_0) = J_1(\theta_0)$

$$\bar{V}_n(\theta) = \frac{1}{n} \sum V_i(\theta), \quad \|\bar{V}_n - v\|_{\infty} \xrightarrow{P} 0$$



$$\begin{aligned} \left\| -\frac{1}{n} \nabla^2 \ell_n(\tilde{\theta}_n) - J_1(\theta_0) \right\| &\leq \left\| \bar{V}_n(\tilde{\theta}_n) - v(\tilde{\theta}_n) \right\| + \left\| v(\tilde{\theta}_n) - v(\theta_0) \right\| \\ &\leq \underbrace{\left\| \bar{V}_n - v \right\|_\infty}_{\xrightarrow{P} 0} + \underbrace{\left\| v(\tilde{\theta}_n) - v(\theta_0) \right\|}_{\xrightarrow{P} 0 \text{ (cts mapping)}} \end{aligned}$$

Hence  $\left( -\frac{1}{n} \nabla^2 \ell_n(\tilde{\theta}_n) \right)^{-1} \xrightarrow{P} J_1(\theta_0)^{-1}$  (cts mapping)

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

Note In this proof we played a bit fast and base with our LLN for random functions, which we only stated for real-valued  $w_n$ .  
So technically we have only justified for  $d=1$ .  
But proof works for  $d>1$ .