

Maximum Likelihood Estimation (MLE)

① Notation

$$1) X \sim f(\cdot | \theta)$$

$$2) L(\theta) = \prod_{i=1}^n f(x_i | \theta) \longrightarrow \text{likelihood}$$

$$\ell(\theta) = \sum_{i=1}^n \ell_i(\theta) = \sum_{i=1}^n \log \{f(x_i | \theta)\}$$

$$3) \hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} L(\theta) = \underset{\theta \in \Theta}{\operatorname{argmax}} \ell(\theta)$$

$\hat{\theta}_{MLE}$ is RV with respect to (x_1, \dots, x_n)

4) Quantity w.r.t log-likelihood $\ell(\theta)$

① Score Function $s(\theta) = \nabla_{\theta} \ell(\theta) \rightarrow \underline{\text{RV}} !!!$

② Fisher Information $I(\theta) = \mathbb{E}_{\theta} [s(\theta) s(\theta)^T]$

$$\begin{aligned}
 & \text{(Differentiability)} \quad \text{Under Some Regularity Condition} \\
 & \left. \begin{array}{l} = -\mathbb{E}_{\theta} [\nabla_{\theta} s(\theta)] \\ = -\mathbb{E}_{\theta} [\nabla_{\theta}^2 \ell(\theta)] = -H(\theta) \end{array} \right\}
 \end{aligned}$$

Note: ① $\mathbb{E}_{\theta} [s(\theta) s(\theta)^T] = \text{Var}_{\theta} [s(\theta)]$ since $\mathbb{E}_{\theta} [s(\theta)] \equiv 0$

→ this quantity will be used in MLE asymptotic property

② $\mathbb{E}_{\theta} [\nabla_{\theta}^2 \ell(\theta)]$ will be used in:

algorithm to solve $\hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} \ell(\theta)$

③ Under some regularity conditions, $\mathbb{E}_{\theta} [s(\theta) s(\theta)^T] = -\mathbb{E}_{\theta} [\nabla_{\theta}^2 \ell(\theta)]$

$$\Rightarrow I(\theta) = -H(\theta)$$

④ Summary:

a) $I(\theta) \rightsquigarrow$ Fisher Information $\rightsquigarrow \mathbb{E}[S(\theta) S(\theta)^T] = \text{Var}[S(\theta)]$

$$\hat{\theta} - \theta_0 \sim N(0, \frac{1}{nI(\theta)})$$

$\Rightarrow I(\theta)$ large \rightarrow variance small \rightarrow efficiency

b) $H(\theta) \rightsquigarrow$ Hessian of log-likelihood $f(\theta)$

$$H(\theta) := \nabla_{\theta}^2 f(\theta)$$



Curvature information!

$\Rightarrow H(\theta)$ large \rightarrow more curved \rightarrow converge faster

② Detailed Discussion

include the Consistency

I. MLE \Leftrightarrow Empirical Risk Minimization (ERM)

Recap: $\hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} f(\theta) = \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_{i=1}^n \log f(x_i | \theta)$

$$\Leftrightarrow \hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \log \frac{f(x_i | \theta_0)}{f(x_i | \theta)}$$

$$= \underset{\theta \in \Theta}{\operatorname{argmin}} R_n(\theta_0, \theta) := \frac{1}{n} \sum_{i=1}^n \log \frac{f(x_i | \theta_0)}{f(x_i | \theta)}$$

Therefore, $\underline{R(\theta_0, \theta)} = \mathbb{E}_{\theta_0} [R_n(\theta_0, \theta)]$

$$= \mathbb{E}_{\theta_0} \left[\log \frac{f(x | \theta_0)}{f(x | \theta)} \right]$$

$$= \int_X \log \frac{f(x|\theta_0)}{f(x|\theta)} \cdot f(x|\theta_0) dx$$

$$= KL(f(x|\theta_0) || f(x|\theta))$$

2 things:

1. From LLN, for fix $\theta \in \Theta$.

$$R_n(\theta_0, \theta) \xrightarrow{P} R(\theta_0, \theta)$$

2. $\theta_0 \in \operatorname{argmin}_{\theta \in \Theta} R(\theta_0, \theta)$

$$= \operatorname{argmin}_{\theta \in \Theta} KL(f(x|\theta_0) || f(x|\theta))$$

Diagram Sketch: $R(\theta_0, \theta) = \mathbb{E}_{\theta_0} [\log \left(\frac{f(X|\theta_0)}{f(X|\theta)} \right)]$

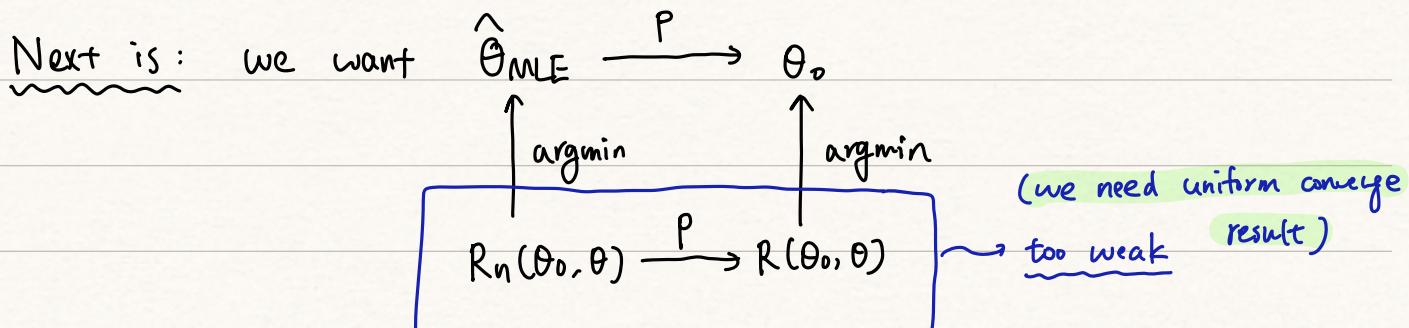
$\mathbb{E}[-\log] \geq -\log \mathbb{E}[]$

$$\begin{aligned} &= -\mathbb{E}_{\theta_0} [\log \frac{f(X|\theta)}{f(X|\theta_0)}] \\ &\geq -\log \mathbb{E}_{\theta_0} \left[\frac{f(X|\theta)}{f(X|\theta_0)} \right] \\ &= -\log \cdot 1 = 0 \end{aligned}$$

To conclude, $\max_{\theta \in \Theta} \ell(\theta) \Leftrightarrow \min_{\theta \in \Theta} R_n(\theta_0, \theta)$



$$\underbrace{R(\theta_0, \theta) := KL(f(x|\theta_0) || f(x|\theta))}_{= \mathbb{E}_{\theta_0} \left[\log \frac{f(x|\theta_0)}{f(x|\theta)} \right]}$$



Thm : Assumption

- a) Strong Identifiability $\rightarrow \inf_{\hat{\theta}: |\hat{\theta} - \theta_0| \geq \varepsilon} KL(p(\theta); p(\hat{\theta})) > 0$
- b) Uniform LLN

"Continuity assumption"

$$\Rightarrow \hat{\theta}_{MLE} \xrightarrow{P} \theta_0 \quad (\text{consistency})$$

Sketch : Strong Identifiability implies that :

if we want $\hat{\theta}_{MLE} \xrightarrow{P} \theta_0$ uniformly

then we only require $KL(f(x|\theta_0) || f(x|\hat{\theta}_{MLE})) \xrightarrow{P} 0$
 $= R(\theta_0, \hat{\theta}_{MLE})$ Sufficiently small !!!

In Binary Classification ERM Framework, we try to analyze:

$$\tilde{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} R(\theta, \theta_0)$$

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} R_n(\theta, \theta_0)$$

$$|R(\hat{\theta}, \theta_0) - R(\tilde{\theta}, \theta_0)| \quad \text{defn of } \hat{\theta}$$

$$= |(R(\hat{\theta}, \theta_0) - R_n(\hat{\theta}, \theta_0)) + (R_n(\hat{\theta}, \theta_0) - R_n(\tilde{\theta}, \theta_0))|$$

$$+ |R_n(\tilde{\theta}, \theta_0) - R(\tilde{\theta}, \theta_0)|$$

Here, actually we assume $\theta_0 \in \Theta$. Thus,

$$R(\theta_0, \hat{\theta}_{MLE}) = R(\theta_0, \hat{\theta}_{MLE}) - R_n(\theta_0, \hat{\theta}_{MLE}) + R_n(\theta_0, \hat{\theta}_{MLE})$$

uniform LLN

$$R_n(\theta_0, \hat{\theta}_{MLE}) \leq R_n(\theta_0, \theta_0) = 0$$

Defn of $\hat{\theta}_{MLE}$

Finish the discussion for
Wasserman LEC 15 Part 4

2. (Extension of 1.)

Suppose $\theta_0 \notin \Theta$

$$(\text{Sketch}) \quad \tilde{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} R(\theta_0, \theta) \quad \underline{KL(f(x|\theta_0) || f(x|\theta))}$$

$$= \underset{\theta \in \Theta}{\operatorname{argmin}} \mathbb{E}_{\theta_0} \left[\log \left(\frac{f(x|\theta_0)}{f(x|\theta)} \right) \right]$$

$$\hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmin}} R_n(\theta_0, \theta)$$

$$= \underset{\theta \in \Theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \log \frac{f(x_i|\theta_0)}{f(x_i|\theta)}$$

Interpretation of $\tilde{\theta}$ $\rightarrow KL(f_{\theta_0} || f_{\tilde{\theta}}) \leq KL(f_{\theta_0} || f_{\theta}) \quad \theta \in \Theta$

\downarrow **
 \tilde{f}_{θ} is the 'KL projection' from f_{θ_0}

to $\{f_\theta : \theta \in \Theta\}$

Similar Analyze $\rightarrow |R(\theta_0, \hat{\theta}_{MLE}) - R(\theta_0, \tilde{\theta})|$

$$= |(R(\theta_0; \hat{\theta}_{MLE}) - R_n(\theta_0; \hat{\theta}_{MLE})) + (R_n(\theta_0; \hat{\theta}_{MLE}) - R_n(\theta_0; \tilde{\theta}))|$$

$\hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{arg\min}} R_n(\theta_0; \theta) + (R_n(\theta_0; \tilde{\theta}) - R(\theta_0; \tilde{\theta}))$

uniform LLN

pointwise LLN

3. Asymptotic Behaviour of $\hat{\theta}_{MLE}$

a) Quantity Recap

score function

$$\textcircled{1} \quad S(X, \theta) = \nabla_\theta f(\theta)$$

$$= \sum_{i=1}^n \nabla_\theta f_i(\theta) = \sum_{i=1}^n \nabla_\theta \log f(X_i | \theta)$$

RV w.r.t (X_1, \dots, X_n)

$S(\theta)$ for simplicity

$$\hookrightarrow S_i(\theta) := \nabla_\theta f_i(\theta)$$

$$\mathbb{E}_\theta [S(\theta)] = \mathbb{E}_\theta [S_i(\theta)] = \mathbb{E}_\theta [\nabla_\theta f_i(\theta)] = 0$$

② Fisher Information Matrix (\underbrace{FIM})

$$(I_1(\theta) := \mathbb{E}_\theta [S_i(\theta) S_i(\theta)^\top] = \text{Var}_\theta [S_i(\theta)]) \rightarrow \text{regular condition (weak)}$$

$$I(\theta) = \mathbb{E}_\theta [S(\theta) S(\theta)^\top] = \mathbb{E}_\theta [-\nabla_\theta^2 f(\theta)] = -H(\theta)$$

$$= \text{Var}_\theta [S(\theta)]$$

$$H_i(\theta) = \mathbb{E}_\theta [\nabla_\theta^2 f_i(\theta)]$$

To conclude,

$$\text{RV } S(\theta) = \sum_{i=1}^n \cdot \nabla_\theta f_i(\theta)$$

expectation 0

variance $I(\theta)$

$$\textcircled{3} \quad \nabla_\theta f(\hat{\theta}_{MLE}) = 0 \quad (\text{roughly speaking})$$

b) **Result** $\sqrt{n}(\hat{\theta}_{MLE} - \theta_0) \sim N(0, I_1(\theta_0)^{-1})$

$$\text{PF: } \nabla_{\theta} l(\hat{\theta}_{MLE}) = \nabla_{\theta} l(\theta_0) + \nabla_{\theta}^2 l(\tilde{\theta})(\hat{\theta}_{MLE} - \theta_0) \\ = 0 \\ \tilde{\theta} = \lambda \theta_0 + (1-\lambda) \hat{\theta}_{MLE} \quad 0 < \lambda < 1$$

$$\Rightarrow \hat{\theta}_{MLE} - \theta_0 = - \nabla_{\theta}^2 l(\tilde{\theta})^{-1} \nabla_{\theta} l(\theta_0) \\ \Leftrightarrow \sqrt{n}(\hat{\theta}_{MLE} - \theta_0) = (-n^{-1} \nabla_{\theta}^2 l(\tilde{\theta}))^{-1} (n^{-\frac{1}{2}} \nabla_{\theta} l(\theta_0))$$

$$\begin{aligned} \textcircled{1} (\text{LLN}) &\Rightarrow -n^{-1} \nabla_{\theta}^2 l(\tilde{\theta}) \xrightarrow{P} -n^{-1} \nabla_{\theta}^2 l(\theta_0) \\ &\xrightarrow{P} -\mathbb{E}_{\theta_0} [\nabla_{\theta}^2 l(\theta_0)] \\ &= I_1(\theta_0) \\ \textcircled{2} (\text{CLT}) &\Rightarrow n^{-\frac{1}{2}} \nabla_{\theta} l(\theta_0) \xrightarrow{d} N(0, \text{var}_{\theta_0}[\nabla_{\theta} l(\theta_0)]) \\ &= N(0, I_1(\theta_0)) \end{aligned}$$

$$\textcircled{3}: \textcircled{1} + \textcircled{2} + \text{Slutsky} \Rightarrow \sqrt{n}(\hat{\theta}_{MLE} - \hat{\theta}_0) \\ = (-n^{-1} \nabla_{\theta}^2 l(\tilde{\theta}))^{-1} (n^{-\frac{1}{2}} \nabla_{\theta} l(\theta_0)) \\ \xrightarrow{d} I_1(\theta_0)^{-1} N(0, I_1(\theta_0)) \\ = N(0, I_1(\theta_0)^{-1})$$

$$\boxed{I_1(\theta_0) := \mathbb{E}_{\theta_0} [S_i(\theta_0) S_i(\theta_0)^T]} \\ = -\mathbb{E}_{\theta_0} [\nabla_{\theta}^2 l(\theta_0)]$$

C) Influence Function Framework

From the previous proof, actually we have shown:

$$\hat{\theta}_{MLE} = \theta_0 - \frac{1}{n} \sum_{i=1}^n \frac{\nabla_{\theta} \log f(x_i | \theta_0)}{I_1(\theta_0)} + \text{reminder}$$

$$= \theta_0 + \frac{1}{n} \sum_{i=1}^n \psi(x_i) \rightarrow \text{more general expression}$$

11/17

GLM Framework Recap

$$\left[\begin{array}{l} \text{Score Function: } S_i(\beta) = \nabla_{\beta} l_i(\beta) \quad \text{RV} \\ \text{Fisher Information Matrix: } I_i(\beta) = \mathbb{E}_{\beta} [S_i(\beta) S_i(\beta)^T] \\ = -\mathbb{E}_{\beta} [\nabla_{\beta}^2 l_i(\beta)] \end{array} \right]$$

$$\text{Observation Matrix} \quad J_i(\beta) = -\nabla_{\beta}^2 l_i(\beta)$$

⇒ Briefly go through

1. Stochastic Part : $Y \sim \text{exponential family}$

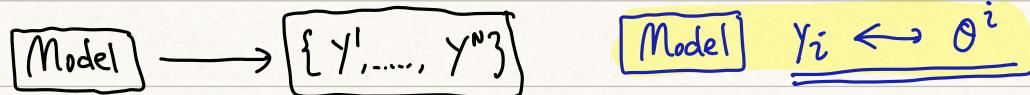
$$\Theta \in \mathbb{R}^q$$

$$f(y|\theta, \phi) = \exp \left\{ \frac{y^T \theta - b(\theta)}{a(\phi)} + c(\phi, y) \right\}$$

Property

$$\textcircled{1} \quad \mathbb{E}_\theta [Y] = \nabla_\theta b(\theta)$$

$$\textcircled{2} \quad \text{Var}_\theta [Y] = a(\phi) \cdot \nabla_\theta^2 b(\theta) \\ := \Sigma(\theta, \phi)$$



$$\mathbb{E}[y^i] = \mu^i = \left[\begin{array}{c} \frac{\partial b}{\partial \theta_1} \\ \vdots \\ \frac{\partial b}{\partial \theta_q} \end{array} \right] \Big|_{\theta=\theta^i}$$

2. Systematic Part : $\eta^i = x^i \cdot \beta$

$$\begin{cases} X^i \in \mathbb{R}^{q \times p} \\ \beta \in \mathbb{R}^p \end{cases}$$

3. Link function : $\underline{g(\mu^i) = \eta^i} \iff \underline{h(\eta^i) = \mu^i}$

$$\begin{bmatrix} \mu^i = \nabla_\theta b(\theta) \end{bmatrix}$$

$$\exists: \eta^i \rightarrow \theta^i$$

General Architecture

$$\theta^i \rightarrow y^i \rightarrow \mathbb{E}_{\theta^i} [Y^i] = \mu^i = \nabla_\theta b(\theta^i) \xleftarrow{\text{Link}} \eta^i = x^i \beta$$

4. MLE $\Rightarrow \hat{\beta}_{MLE} \Rightarrow \arg \max \ell(\beta)$

$$\ell(\beta) = \sum_{i=1}^N \ell_i(\beta)$$

$$= \sum_{i=1}^N \left\{ \frac{(y_i)^T Q^i - b(\theta^i)}{a(\phi)} + \text{const} \right\}$$

denote as $\ell_i(\beta)$

$$\text{Then, } \hat{\beta}_{MLE} = \underset{\beta}{\operatorname{argmax}} \sum_{i=1}^N \ell_i(\beta)$$

5. One algo is, (Fisher Scoring algo)

Newton Method

$$\rightarrow \hat{\beta}_{t+1} = \hat{\beta}_t + I(\hat{\beta}_t)^{-1} S(\hat{\beta}_t)$$

Remark:

→ in the conventional frameworks

$$\hat{\beta}_{t+1} \leftarrow \hat{\beta}_t + J(\hat{\beta}_t)^{-1} S(\hat{\beta}_t)$$

→ Here, replace $J(\hat{\beta}_t)$ with $I(\hat{\beta}_t)$

Tricks

$$S(\beta) = \nabla_{\beta} \ell(\beta)$$

$$I(\beta) = - \mathbb{E}_{\beta} [\nabla_{\beta}^2 \ell(\beta)] = \mathbb{E}_{\beta} [S(\beta) S(\beta)^T]$$

$$J(\beta) = - \nabla_{\beta}^2 \ell(\beta)$$

$$\text{Define } S_i(\beta) = \nabla_{\beta} \ell_i(\beta) \quad I_i(\beta) = - \mathbb{E}_{\beta} [\nabla_{\beta}^2 \ell_i(\beta)] = \mathbb{E}_{\beta} [J_i(\beta)]$$

$$\text{Then } J_i(\beta) = - \nabla_{\beta}^2 \ell_i(\beta) = I_i(\beta) - \frac{1}{a(\phi)} R_i(\beta)$$

$\downarrow \text{fix}$ $\downarrow \text{random part}$

Note: When applying Canonical Link Function : $J_i(\beta) = \mathbb{E}_{\beta} [J_i(\beta)] = I_i(\beta)$

→ that is to say, $J_i(\beta)$ is not a RV!

After Careful Calculation, the Fisher Scoring Algo

turns to

Iterative Reweighted Least Square

- { for $\mathbf{y} \sim 1\text{-dim}$ → check Lecture 3 in ST5213
- for $\mathbf{y} \sim g\text{-dim}$ → check GLM in ST5213

6. Asymptotic Result for $\hat{\beta}_{MLE}$

$$\textcircled{1} \quad \begin{cases} \hat{\beta}_{MLE} - \beta_0 \xrightarrow{d} N(0, I_1(\beta_0)^{-1}) \\ \hat{\beta}_{MLE} - \beta_0 \xrightarrow{d} N(0, J_1(\beta_0)^{-1}) \end{cases}$$

② Hypothesis Testing

a) LRT $\left\{ \begin{array}{l} H_0: \beta = \text{given vector} \\ H_A: \text{otherwise} \end{array} \right.$

$$-2 \log \Lambda \xrightarrow{d} \chi^2_{df}$$

b) Wald Testing ← asymptotic normality

$$I(\beta) = \mathbb{E}_\beta [S(\beta) S(\beta)^T]$$

closed form

standard error

- { ① exact $SE(\hat{\beta}_j) = \text{diag}(I(\hat{\beta}))_j$
- ② Bootstrap to estimate SE if the sample size is moderate

Asymptotic Variance may have large BIAS