

# Maximum Likelihood Estimates: Asymptotic Properties

Ching-Kang Ing

Institute of Statistics, National Tsing Hua University, Hsinchu, Taiwan

# Outline

- 1 Consistency of MLEs
- 2 Logistic regression model
- 3 A central limit theorem for MLEs
  - Example
- 4 Asymptotics for likelihood ratio tests
- 5 Pearson's Chi-Squared Test with unknown parameters
- 6 Akaike's information criterion (AIC)
- 7 More on information criteria
  - Bayesian information criterion (BIC)
- 8 A Newton-Raphson method and its asymptotics
- 9 Asymptotics for the MLEs of Weibull distribution parameters

# Consistency of MLEs

Let  $L(\boldsymbol{\theta})$  be a likelihood function and  $\boldsymbol{\theta}_0 \in \boldsymbol{\Theta}$  denote the true parameter. If

- (i)  $\boldsymbol{\theta}_0$  is an interior point of  $\boldsymbol{\Theta}$ ,
- (ii) There exists a sufficiently small number  $a > 0$  and a nonrandom matrix  $\mathbf{G}(\boldsymbol{\theta})$  such that for any  $\boldsymbol{\theta} \in \bar{B}_a(\boldsymbol{\theta}_0) = \{\boldsymbol{\theta} : \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| \leq a\}$ ,

$$\sup_{\boldsymbol{\theta} \in \bar{B}_a(\boldsymbol{\theta}_0)} \left\| -\frac{1}{n} \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} - \mathbf{G}(\boldsymbol{\theta}) \right\| \rightarrow 0 \quad \text{in probability,}$$

$$\inf_{\boldsymbol{\theta} \in \bar{B}_a(\boldsymbol{\theta}_0)} \lambda_{\min}(\mathbf{G}(\boldsymbol{\theta})) > 0,$$

where  $\ell(\boldsymbol{\theta}) = \log L(\boldsymbol{\theta})$ ,

- (iii)  $\frac{\partial}{\partial \boldsymbol{\theta}} \ell(\boldsymbol{\theta}_0) = O_p(\sqrt{n})$ ,

## Norm

- For a vector  $\mathbf{v}$ ,  $\|\mathbf{v}\|$  denotes the Euclidean norm of  $\mathbf{v}$ .
- For a matrix  $\mathbf{V}$ ,  $\|\mathbf{V}\|$  denotes the spectral norm of  $\mathbf{V}$ .

then the likelihood equation

$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbf{0} \quad (1)$$

has a root  $\tilde{\boldsymbol{\theta}}_n$  satisfying

$$\tilde{\boldsymbol{\theta}}_n \rightarrow \boldsymbol{\theta}_0 \quad \text{in probability.} \quad (2)$$

Moreover, if

$$(iv) \quad P\left(\ell(\boldsymbol{\theta}_0) > \sup_{\boldsymbol{\theta} \in \bar{B}_a^c(\boldsymbol{\theta}_0)} \ell(\boldsymbol{\theta})\right) \rightarrow 1,$$

then

$$\hat{\boldsymbol{\theta}}_n = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta}) \rightarrow \boldsymbol{\theta}_0 \quad \text{in probability.} \quad (3)$$

# Proof

- Let  $\tilde{\theta}_n$  be the solution of (1) closest to  $\theta_0$ . We will first show that for any  $0 < \epsilon \leq a$  (note that  $a$  is defined in (ii)),

$$P\left(\sup_{\theta \in B_\epsilon^*(\theta_0)} \ell(\theta) < \ell(\theta_0)\right) \rightarrow 1 \text{ as } n \rightarrow \infty, \quad (4)$$

where  $B_\epsilon^*(\theta_0) = \bar{B}_a(\theta_0) - B_\epsilon(\theta_0)$ .

- Relation (4) implies

$$P\left(\tilde{\theta}_n \in B_\epsilon(\theta_0)\right) = P(\text{there exists a solution of (1) lying in } B_\epsilon(\theta_0)) \rightarrow 1 \text{ as } n \rightarrow \infty,$$

and hence the desired conclusion (2) follows.

- Define

$$E_n(\delta) = \left\{ \inf_{\theta \in \bar{B}_a(\theta_0)} \lambda_{\min} \left( -\frac{1}{n} \frac{\partial^2 \ell(\theta)}{\partial \theta_i \partial \theta_j} \right) < \delta \right\},$$

where  $\delta > 0$  is arbitrarily small.

Then, on the set  $E_n^c(\delta)$  (the complement of  $E_n(\delta)$ ), one has by Taylor's theorem,

$$\begin{aligned}
 & \sup_{\boldsymbol{\theta} \in B_\epsilon^*(\boldsymbol{\theta}_0)} \ell(\boldsymbol{\theta}) - \ell(\boldsymbol{\theta}_0) \\
 & \stackrel{\text{why?}}{\leq} \sup_{\boldsymbol{\theta} \in B_\epsilon^*(\boldsymbol{\theta}_0)} \left\| \frac{\partial \ell(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \right\| \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| \\
 & \quad - \inf_{\boldsymbol{\theta} \in B_\epsilon^*(\boldsymbol{\theta}_0)} \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)' \left( -\frac{\partial^2 \ell(\boldsymbol{\theta}^*)}{\partial \theta_i \partial \theta_j} \right) (\boldsymbol{\theta} - \boldsymbol{\theta}_0) \quad (\|\boldsymbol{\theta}^* - \boldsymbol{\theta}_0\| \leq \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|) \\
 & \leq \left\| \frac{\partial \ell(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \right\| a - \frac{\epsilon^2}{2} \inf_{\boldsymbol{\theta} \in \bar{B}_a(\boldsymbol{\theta}_0)} \lambda_{\min} \left( -\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} \right) \\
 & \leq \left\| \frac{\partial \ell(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \right\| a - \frac{n\epsilon^2}{2} \delta.
 \end{aligned}$$

This and condition (iii) yield

$$P \left( \sup_{\boldsymbol{\theta} \in B_\epsilon^*(\boldsymbol{\theta}_0)} \ell(\boldsymbol{\theta}) - \ell(\boldsymbol{\theta}_0) \geq 0, E_n^c(\delta) \right) \leq P \left( \left\| \frac{\partial \ell(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \right\| a \geq \frac{\epsilon^2}{2} n\delta \right) \rightarrow 0 \text{ as } n \rightarrow \infty. \quad (5)$$

- In view of (5), (4) follows if we can show that for any  $\epsilon_1 > 0$ , there exists a sufficiently large  $N$  and sufficiently small  $\delta > 0$  such that for all  $n \geq N$ ,

$$P(E_n(\delta)) < \epsilon_1. \quad (6)$$

- To show (6), denote  $-\frac{1}{n} \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j}$  by  $\mathbf{A}_n(\boldsymbol{\theta})$ . Then,

$$\begin{aligned} \lambda_{\min}(\mathbf{A}_n(\boldsymbol{\theta})) &\geq \lambda_{\min}(\mathbf{A}_n(\boldsymbol{\theta}) - \mathbf{G}(\boldsymbol{\theta})) + \lambda_{\min}(\mathbf{G}(\boldsymbol{\theta})) \\ &\geq \lambda_{\min}(\mathbf{G}(\boldsymbol{\theta})) - \|\mathbf{A}_n(\boldsymbol{\theta}) - \mathbf{G}(\boldsymbol{\theta})\|, \end{aligned}$$

yielding

$$\inf_{\boldsymbol{\theta} \in \bar{B}_a(\boldsymbol{\theta}_0)} \lambda_{\min}(\mathbf{A}_n(\boldsymbol{\theta})) \geq \inf_{\boldsymbol{\theta} \in \bar{B}_a(\boldsymbol{\theta}_0)} \lambda_{\min}(\mathbf{G}(\boldsymbol{\theta})) - \sup_{\boldsymbol{\theta} \in \bar{B}_a(\boldsymbol{\theta}_0)} \|\mathbf{A}_n(\boldsymbol{\theta}) - \mathbf{G}(\boldsymbol{\theta})\|.$$

This and condition (ii) give (6). Thus, (2) is proved.

- Finally, it follows from (4) and (iv) that

$$P\left(\sup_{\boldsymbol{\theta} \in B_\epsilon^c(\boldsymbol{\theta}_0)} \ell(\boldsymbol{\theta}) < \ell(\boldsymbol{\theta}_0)\right) \rightarrow 1 \text{ as } n \rightarrow \infty,$$

which implies  $\hat{\boldsymbol{\theta}}_n$  converges  $\boldsymbol{\theta}_0$  in probability.

# Logistic Regression Model

- Recall logistic regression. For  $i = 1, \dots, n$ , assume

$$P(Y_i = 1) = 1 - P(Y_i = 0) = \frac{e^{\mathbf{x}'_i \beta_0}}{1 + e^{\mathbf{x}'_i \beta_0}} \equiv p_{i,0}, \quad \left( \log \left( \frac{p_{i,0}}{1 - p_{i,0}} \right) = \mathbf{x}'_i \beta_0 \right)$$

where  $\mathbf{x}_i$  is the explanatory vector and  $\beta_0$  is an unknown regression coefficient vector.

- Therefore, the likelihood function is (assuming  $Y_i$  are independent)

$$L(\beta) = P(Y_1 = y_1, \dots, Y_n = y_n) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i},$$

where  $p_i = e^{\mathbf{x}'_i \beta} / (1 + e^{\mathbf{x}'_i \beta})$  and  $\beta \in B$ , the parameter space, and the log-likelihood function is

$$\begin{aligned} \ell(\beta) &= \log L(\beta) = \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)] \\ &= \sum_{i=1}^n y_i \log \left( \frac{p_i}{1 - p_i} \right) + \sum_{i=1}^n \log(1 - p_i) = \sum_{i=1}^n y_i \mathbf{x}'_i \beta - \sum_{i=1}^n \log(1 + e^{\mathbf{x}'_i \beta}). \end{aligned}$$

- Moreover,

$$\frac{\partial \ell(\beta)}{\partial \beta} = \sum_{i=1}^n \mathbf{x}_i y_i - \sum_{i=1}^n \frac{e^{\mathbf{x}'_i \beta}}{1 + e^{\mathbf{x}'_i \beta}} \mathbf{x}_i = \sum_{i=1}^n \mathbf{x}_i \left( y_i - \frac{e^{\mathbf{x}'_i \beta}}{1 + e^{\mathbf{x}'_i \beta}} \right) \Big|_{\beta = \beta_0} = \sum_{i=1}^n \mathbf{x}_i (y_i - p_{i,0}).$$

- Note that  $y_i = E(y_i) + (y_i - E(y_i)) = p_{i,0} + \epsilon_i$  where  $\epsilon_i \sim (0, p_{i,0}(1 - p_{i,0}))$ .



Define  $\hat{\beta} = \operatorname{argmax}_{\beta \in B} \ell(\beta)$ . Does  $\hat{\beta}$  converges to  $\beta_0$  in probability?

(0)  $\beta_0$  is an interior point of  $B$ , ← This can be assumed.

$$(1) \frac{1}{\sqrt{n}} \frac{\partial \ell(\beta_0)}{\partial \beta} = O_p(1),$$

$$\frac{1}{\sqrt{n}} \frac{\partial \ell(\beta_0)}{\partial \beta} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i \epsilon_i \text{ in linear regression model.}$$

(2)  $P(\inf_{\beta \in \overline{B}_a(\beta_0)} \lambda_{\min}(-\frac{1}{n} \frac{\partial^2 \ell(\beta)}{\partial \beta_i \partial \beta_j}) > \delta) \rightarrow 1$ , where  $\delta$  and  $a$  are small constants and

$$\overline{B}_a(\beta_0) = \{\beta : \|\beta - \beta_0\| \leq a\},$$

$$-\frac{1}{n} \frac{\partial^2 \ell(\beta)}{\partial \beta_i \partial \beta_j} = \frac{1}{n} \mathbf{X}' \mathbf{X} \text{ in linear regression model.}$$

(3)  $P(\ell(\beta_0) > \sup_{\beta \in \overline{B}_a^c(\beta_0)} \ell(\beta)) \rightarrow 1$ .

In linear regression model,

$$-\frac{1}{2\sigma^2} \sum_{i=1}^n \epsilon_i^2 = -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i' \beta_0)^2 > -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i' \beta)^2, \quad \beta \in \overline{B}_a^c(\beta_0),$$

because

$$-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i' \beta)^2 = -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i' \beta_0 + \mathbf{x}_i' (\beta_0 - \beta))^2$$

$$= -\frac{1}{2\sigma^2} \sum_{i=1}^n (\epsilon_i + \mathbf{x}_i' (\beta_0 - \beta))^2$$

$$\sim -\frac{1}{2\sigma^2} \left\{ \sum_{i=1}^n \epsilon_i^2 + (\beta_0 - \beta)' \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right) (\beta_0 - \beta) \right\}.$$

## Remark

- For (1), does  $\frac{1}{\sqrt{n}} \frac{\partial \ell(\beta_0)}{\partial \beta} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i (y_i - p_{i,0}) = O_p(1)$ ? Note that

$$\begin{aligned}
 \frac{1}{n} E \left( \frac{\partial \ell(\beta_0)}{\partial \beta} \left( \frac{\partial \ell(\beta_0)}{\partial \beta} \right)' \right) &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \text{Var}(y_i) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' p_{i,0} (1 - p_{i,0}) \\
 &= \frac{1}{n} \mathbf{X}' D_0 \mathbf{X} \\
 &\leq \lambda_{\max} \left( \frac{1}{n} \mathbf{X}' \mathbf{X} \right) \lambda_{\max}(D_0) \\
 &\leq \lambda_{\max} \left( \frac{1}{n} \mathbf{X}' \mathbf{X} \right) \min_{1 \leq i \leq n} p_{i,0} (1 - p_{i,0}) \\
 &\leq \frac{1}{4} \lambda_{\max} \left( \frac{1}{n} \mathbf{X}' \mathbf{X} \right).
 \end{aligned}$$

where  $D_0 = \text{diag}(p_{1,0}(1 - p_{1,0}), \dots, p_{n,0}(1 - p_{n,0}))$ . If  $\lambda_{\max}(\frac{1}{n} \mathbf{X}' \mathbf{X})$  is bounded, then (1) holds.

- For (2), since  $-\frac{1}{n} \frac{\partial^2 \ell(\beta)}{\partial \beta_i \partial \beta_j} \sim \frac{1}{n} \mathbf{X}' D \mathbf{X}$  where  $D = \text{diag}(p_1(1 - p_1), \dots, p_n(1 - p_n))$ , we have

$$\inf_{\beta \in \overline{B}_a(\beta_0)} \lambda_{\min} \left( \frac{1}{n} \mathbf{X}' D \mathbf{X} \right) \geq \lambda_{\min} \left( \frac{1}{n} \mathbf{X}' \mathbf{X} \right) \inf_{\beta \in \overline{B}_a(\beta_0)} \min_{1 \leq i \leq n} p_i (1 - p_i).$$

Therefore, (2) holds if  $\lambda_{\min}(\frac{1}{n} \mathbf{X}' \mathbf{X}) > \underline{c}_1 > 0$  and  $\inf_{\beta \in B} \min_{1 \leq i \leq n} p_i (1 - p_i) > \underline{c}_2 > 0$  for all large  $n$ , where  $\underline{c}_1$  and  $\underline{c}_2$  are some positive constants.

## Remark (cont.)

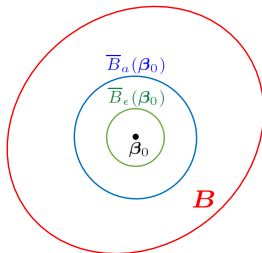
- For (3), it can be shown through (0)–(2) that

$$P \left( \ell(\beta_0) > \sup_{\beta \in (\overline{B}_a(\beta_0) - B_\epsilon(\beta_0))} \ell(\beta) \right) \rightarrow 1,$$

for any  $\epsilon > 0$ . Moreover, since

$$\frac{\partial^2 \ell(\beta)}{\partial \beta_i \partial \beta_j} \sim -\mathbf{X}' \mathbf{D} \mathbf{X},$$

which is negative definite (assuming  $\mathbf{X}' \mathbf{X}$  is p.d. and  $\min_{1 \leq i \leq n} p_i(1 - p_i)$  is bounded away from 0 in  $B$ ) and hence  $\ell(\beta)$  is a convex function, (3) follows.



# A central limit theorem for MLEs

- Define

$$\hat{\eta} = \underset{\eta \in \Lambda}{\operatorname{argmin}} \ell(\eta),$$

where  $\ell(\eta)$  is the log-likelihood function and the true parameter  $\eta_0$  is an interior point of the parameter space  $\Lambda$ .

- I have shown in my previous note that

$$\hat{\eta} \xrightarrow{pr.} \eta_0.$$

- In the following, I will present a CLT for  $\hat{\eta}$  using a way far from being rigorous.

- Note that

$$\mathbf{0} = \frac{\partial \ell(\hat{\boldsymbol{\eta}})}{\partial \boldsymbol{\eta}} = \frac{\partial \ell(\boldsymbol{\eta}_0)}{\partial \boldsymbol{\eta}} + \frac{\partial^2 \ell(\boldsymbol{\eta}^*)}{\partial \eta_i \partial \eta_j} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0), \quad (7)$$

where  $\|\boldsymbol{\eta}^* - \boldsymbol{\eta}_0\| \leq \|\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0\|$ .

- Since it can be shown that  $\boldsymbol{\eta}^* \rightarrow \boldsymbol{\eta}_0$  in probability (because  $\hat{\boldsymbol{\eta}} \rightarrow \boldsymbol{\eta}_0$  in probability), we have by (7) and some tedious arguments (which we skip in this note),

$$\sqrt{n}(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0) - \left[ - \left( \frac{1}{n} \frac{\partial^2 \ell(\boldsymbol{\eta}_0)}{\partial \eta_i \partial \eta_j} \right)^{-1} \left( \frac{1}{\sqrt{n}} \frac{\partial \ell(\boldsymbol{\eta}_0)}{\partial \boldsymbol{\eta}} \right) \right] \xrightarrow{pr.} \mathbf{0}. \quad (8)$$

( $\boldsymbol{\eta}^*$  in (7) has been replaced here.)

Assume

$$E \left( -\frac{1}{n} \frac{\partial^2 \ell(\boldsymbol{\eta}_0)}{\partial \eta_i \partial \eta_j} \right) = E \left[ \frac{1}{n} \left( \frac{\partial \ell(\boldsymbol{\eta}_0)}{\partial \boldsymbol{\eta}} \right) \left( \frac{\partial \ell(\boldsymbol{\eta}_0)}{\partial \boldsymbol{\eta}} \right)' \right] \xrightarrow[n \rightarrow \infty]{} \mathbf{I}(\boldsymbol{\eta}_0),$$

and

$$\frac{1}{n} \frac{\partial^2 \ell(\boldsymbol{\eta}_0)}{\partial \eta_i \partial \eta_j} - E \left( \frac{1}{n} \frac{\partial^2 \ell(\boldsymbol{\eta}_0)}{\partial \eta_i \partial \eta_j} \right) \xrightarrow{pr.} \mathbf{0}. \text{ (zero matrix)}$$

Then

$$\frac{1}{n} \frac{\partial^2 \ell(\boldsymbol{\eta}_0)}{\partial \eta_i \partial \eta_j} \xrightarrow{pr.} -\mathbf{I}(\boldsymbol{\eta}_0),$$

and hence

$$-\left( \frac{1}{n} \frac{\partial^2 \ell(\boldsymbol{\eta}_0)}{\partial \eta_i \partial \eta_j} \right)^{-1} \xrightarrow{pr.} \mathbf{I}^{-1}(\boldsymbol{\eta}_0). \quad (9)$$

- Under certain regularity conditions, it can be shown that

$$\frac{1}{\sqrt{n}} \frac{\partial \ell(\boldsymbol{\eta}_0)}{\partial \boldsymbol{\eta}} \xrightarrow{d} N(\mathbf{0}, \mathbf{I}(\boldsymbol{\eta}_0)). \quad (10)$$

- By (8)–(10) and Slutsky's Theorem,

$$\sqrt{n}(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{I}^{-1}(\boldsymbol{\eta}_0)). \quad (11)$$

# Example I

Consider logistic regression model. Does  $\sqrt{n}(\hat{\beta} - \beta_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{A})$ , for some  $\mathbf{A}$ ?

- By  $\frac{\partial \ell(\hat{\beta})}{\partial \beta} \sim \frac{\partial \ell(\beta_0)}{\partial \beta} + \frac{\partial^2 \ell(\beta_0)}{\partial \beta_i \partial \beta_j} (\hat{\beta} - \beta_0)$ , we have

$$\sqrt{n}(\hat{\beta} - \beta_0) \sim \left( -\frac{1}{n} \frac{\partial^2 \ell(\beta_0)}{\partial \beta_i \partial \beta_j} \right)^{-1} \frac{1}{\sqrt{n}} \frac{\partial \ell(\beta_0)}{\partial \beta}.$$

- Assume  $\frac{1}{n} \frac{\partial^2 \ell(\beta_0)}{\partial \beta_i \partial \beta_j} - E \left( \frac{1}{n} \frac{\partial^2 \ell(\beta_0)}{\partial \beta_i \partial \beta_j} \right) \xrightarrow{pr.} 0$  and

$$E \left( -\frac{1}{n} \frac{\partial^2 \ell(\beta_0)}{\partial \beta_i \partial \beta_j} \right) = E \left( \frac{1}{n} \frac{\partial \ell(\beta_0)}{\partial \beta} \left( \frac{\partial \ell(\beta_0)}{\partial \beta} \right)' \right) \xrightarrow{n \rightarrow \infty} \mathbf{G} \text{ (p.d.)}.$$

Then

$$-\frac{1}{n} \frac{\partial^2 \ell(\beta_0)}{\partial \beta_i \partial \beta_j} \sim \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' p_{i,0} (1 - p_{i,0}) \xrightarrow{n \rightarrow \infty} \mathbf{G}.$$

- Since  $\frac{1}{\sqrt{n}} \frac{\partial \ell(\beta_0)}{\partial \beta} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i \epsilon_i$ , it follows from Lindeberg's CLT that

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i \epsilon_i \xrightarrow{d} N(\mathbf{0}, \mathbf{G}),$$

and hence  $\sqrt{n}(\hat{\beta} - \beta_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{G}^{-1})$ .



## Example II

- Consider  $y_i \overset{\text{independent}}{\sim} N(\mathbf{x}'_i \boldsymbol{\beta}_0, e^{\mathbf{x}'_i \boldsymbol{\alpha}_0})$ , or equivalently,

$$y_i = \mathbf{x}'_i \boldsymbol{\beta}_0 + \sigma_{i,0} \epsilon_i,$$

with  $\sigma_{i,0} = e^{\frac{1}{2} \mathbf{x}'_i \boldsymbol{\alpha}_0}$  and  $\epsilon_i \overset{\text{i.i.d.}}{\sim} N(0, 1)$ .

- log-likelihood function

$$\ell(\boldsymbol{\eta}) = \ell(\boldsymbol{\beta}, \boldsymbol{\alpha}) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^n \mathbf{x}'_i \boldsymbol{\alpha} - \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}'_i \boldsymbol{\beta})^2 e^{-\mathbf{x}'_i \boldsymbol{\alpha}}.$$

- To illustrate (11), we start by considering the special case where  $\sigma_{i,0} = 1$  or equivalently,  $\alpha_0 = 0$ . Then,

$$y_i = \mathbf{x}'_i \beta_0 + \epsilon_i, \quad \epsilon_i \stackrel{\text{i.i.d.}}{\sim} N(0, 1).$$

- In addition,

$$\ell(\beta) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}'_i \beta)^2,$$

yielding

$$\frac{\partial \ell(\beta_0)}{\partial \beta} = \sum_{i=1}^n \mathbf{x}_i (y_i - \mathbf{x}'_i \beta_0) = \sum_{i=1}^n \mathbf{x}_i \epsilon_i.$$

- Note that we've shown before

$$\frac{1}{\sqrt{n}} \frac{\partial \ell(\beta_0)}{\partial \beta} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i \epsilon_i \xrightarrow{d} N(\mathbf{0}, \mathbf{R}),$$

where  $\mathbf{R} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}'_i$  (assuming the limit exists).

- Moreover, we have

$$\frac{\partial^2 \ell(\beta_0)}{\partial \beta_i \partial \beta_j} = - \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i',$$

hence

$$\frac{1}{n} \frac{\partial^2 \ell(\beta_0)}{\partial \beta_i \partial \beta_j} \rightarrow -\mathbf{R}.$$

- As a result, we have by (11),

$$\sqrt{n}(\hat{\beta} - \beta_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{R}^{-1}),$$

which coincides with the result that we've obtained for LSE. (Noting that in this example, we assume  $\sigma^2 = 1$ .)

We now get back to the general case  $\alpha_0 \neq 0$ .

$$\frac{\partial \ell(\eta_0)}{\partial \beta} = \sum_{i=1}^n \frac{\mathbf{x}_i}{\sigma_{i,0}} \epsilon_i,$$

$$\frac{\partial \ell(\eta_0)}{\partial \alpha} = \frac{1}{2} \sum_{i=1}^n \left( \frac{(y_i - \mathbf{x}_i' \beta_0)^2}{\sigma_{i,0}^2} - 1 \right) \mathbf{x}_i, \quad \left( \frac{(y_i - \mathbf{x}_i' \beta_0)^2}{\sigma_{i,0}^2} = \epsilon_i^2 \right)$$

$$\frac{\partial^2 \ell(\eta_0)}{\partial \beta_i \partial \beta_j} = - \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i'}{\sigma_{i,0}^2},$$

$$\frac{\partial^2 \ell(\eta_0)}{\partial \alpha_i \partial \alpha_j} = - \frac{1}{2} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \epsilon_i^2,$$

$$\frac{\partial^2 \ell(\eta_0)}{\partial \alpha_i \partial \beta_j} = - \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i'}{\sigma_{i,0}^2} (y_i - \mathbf{x}_i' \beta_0).$$

- Assume

$$\frac{1}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i'}{\sigma_{i,0}^2} \xrightarrow{n \rightarrow \infty} \mathbf{B} \quad \text{and} \quad \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \xrightarrow{n \rightarrow \infty} \mathbf{R}.$$

Then,

$$-\frac{1}{n} \frac{\partial^2 \ell(\boldsymbol{\eta}_0)}{\partial \eta_i \partial \eta_j} \xrightarrow{pr.} \begin{pmatrix} \mathbf{B} & \mathbf{0} \\ \mathbf{0} & \frac{1}{2} \mathbf{R} \end{pmatrix}.$$

- Moreover, it can be shown that

$$\frac{1}{\sqrt{n}} \frac{\partial \ell(\boldsymbol{\eta}_0)}{\partial \boldsymbol{\eta}} = \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\mathbf{x}_i}{\sigma_{i,0}} \epsilon_i, \frac{1}{2\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i (\epsilon_i^2 - 1) \right) \xrightarrow{d} N \left( \mathbf{0}, \begin{pmatrix} \mathbf{B} & \mathbf{0} \\ \mathbf{0} & \frac{1}{2} \mathbf{R} \end{pmatrix} \right).$$

Consequently,

$$\sqrt{n}(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0) \xrightarrow{d} N \left( \mathbf{0}, \begin{pmatrix} \mathbf{B}^{-1} & \mathbf{0} \\ \mathbf{0} & 2\mathbf{R}^{-1} \end{pmatrix} \right).$$

### Question

If  $\epsilon_i \stackrel{\text{i.i.d.}}{\sim} (0, 1)$  (without assuming normality;  $E(\epsilon_i) = 0$  and  $Var(\epsilon_i) = 1$ ),  $\sqrt{n}(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0) \xrightarrow{d} ?$

- In practice,  $B$  and  $R$  are unknown and we can use

$$\hat{B} = \frac{1}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i'}{e^{\mathbf{x}_i' \hat{\alpha}}} \quad \text{and} \quad \hat{R} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i'$$

in place of  $B$  and  $R$ .

- Moreover, it can be shown that

$$\sqrt{n} \begin{pmatrix} \hat{B}^{\frac{1}{2}} & \mathbf{0} \\ \mathbf{0} & \frac{1}{\sqrt{2}} \hat{R}^{\frac{1}{2}} \end{pmatrix} (\hat{\eta} - \eta_0) \xrightarrow{d} N(\mathbf{0}, I),$$

which allow us to construct confidence regions for  $\eta_0$ .

### Question

If  $\epsilon_i \stackrel{\text{i.i.d.}}{\sim} (0, 1)$  (without assuming normality;  $E(\epsilon_i) = 0$  and  $Var(\epsilon_i) = 1$ ), please find

a "data-driven matrix"  $\hat{D}$

such that

$$\sqrt{n} \hat{D} (\hat{\eta} - \eta_0) \xrightarrow{d} N(\mathbf{0}, I).$$

# Asymptotics for likelihood ratio tests

- Let  $x_1, \dots, x_n$  be a random sample from the pdf  $f_{\theta_0}(\cdot)$ , where

$$\theta_0 = \begin{pmatrix} \theta_{0,K(r)} \\ \theta_{0,r} \end{pmatrix} \in \Theta \subseteq \mathbb{R}^K,$$

$\theta_0$  is an interior point of  $\Theta$ ,  $\theta_{0,r}$  and  $\theta_{0,K(r)}$ , respectively, are  $r$ -dimensional and  $(K - r)$ -dimensional vectors, and  $K > r$ .

- For the hypothesis

$$H_0 : \theta_0 \in \Theta_0 \subseteq \Theta \quad \text{versus} \quad H_A : \theta_0 \notin \Theta_0,$$

the likelihood ratio test rejects  $H_0$  if  $\Lambda_n \leq C$ , where

$$\Lambda_n = \frac{\sup_{\theta \in \Theta_0} L(\theta)}{\sup_{\theta \in \Theta} L(\theta)} = \frac{\sup_{\theta \in \Theta_0} \prod_{i=1}^n f_{\theta}(x_i)}{\sup_{\theta \in \Theta} \prod_{i=1}^n f_{\theta}(x_i)},$$

and  $0 \leq C \leq 1$  is determined by the distribution of  $\Lambda_n$  and the level of the test.

- In the following, I'll show that

$$-2 \log \Lambda_n \xrightarrow{d} \chi^2(\dim(\Theta) - \dim(\Theta_0)).$$

- Without loss of generality, we may assume

$$H_0 : \theta_{0,r} = \theta_{0,r}^* \quad \text{and} \quad H_A : \theta_{0,r} \neq \theta_{0,r}^*,$$

and hence

$$\Theta_0 = \{\theta : (\theta'_{K(r)}, \theta_{0,r}^{*'}) \in \Theta\},$$

where  $\theta_{K(r)}$  is  $(K - r)$ -dimensional.

( $\theta_{K(r)}$ : free parameters;  $\Theta$ : non-degenerate subset in  $\mathbb{R}^K$ )



Assume

$$\hat{\theta}^* = \begin{pmatrix} \hat{\theta}_{K(r)} \\ \theta_{0,r}^* \end{pmatrix} = \underset{\theta \in \Theta_0}{\operatorname{argmax}} L(\theta) \quad \text{and} \quad \hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmax}} L(\theta)$$

are unique. Then, under some regularity conditions mentioned in my previous note,

$$\hat{\theta} = \begin{pmatrix} \hat{\theta}_1 \\ \vdots \\ \hat{\theta}_K \end{pmatrix} \xrightarrow{pr.} \theta_0 \equiv \begin{pmatrix} \theta_{0,1} \\ \vdots \\ \theta_{0,K-r} \\ \theta_{0,K-r+1} \\ \vdots \\ \theta_{0,K} \end{pmatrix} = \begin{pmatrix} \theta_{0,K(r)} \\ \theta_{0,r} \end{pmatrix}, \quad \hat{\theta} \xrightarrow[H_0]{pr.} \theta_0^* \equiv \begin{pmatrix} \theta_{0,K(r)} \\ \theta_{0,r}^* \end{pmatrix} = \begin{pmatrix} \theta_{0,1} \\ \vdots \\ \theta_{0,K-r} \\ \theta_{0,K-r+1}^* \\ \vdots \\ \theta_{0,K}^* \end{pmatrix},$$

and

$$\hat{\theta}^* = \begin{pmatrix} \hat{\theta}_{K(r)} \\ \theta_{0,r}^* \end{pmatrix} = \begin{pmatrix} \hat{\theta}_{K(r),1} \\ \vdots \\ \hat{\theta}_{K(r),K-r} \\ \theta_{0,K-r+1}^* \\ \vdots \\ \theta_{0,K}^* \end{pmatrix} \xrightarrow[H_0]{pr.} \theta_0^*.$$

Now, by Taylor's theorem, we have under  $H_0$ ,

$$\begin{aligned}\ell(\hat{\theta}^*) &= \ell(\theta_0^*) + \left( \frac{\partial \ell(\theta_0^*)}{\partial \theta_{K(r)}} \right)' (\hat{\theta}_{K(r)} - \theta_{0,K(r)}) \\ &\quad + \frac{1}{2} (\hat{\theta}_{K(r)} - \theta_{0,K(r)})' \left( \frac{\partial^2 \ell(\theta_0^*)}{\partial \theta_i \partial \theta_j} \right)_{1 \leq i, j \leq K-r} (\hat{\theta}_{K(r)} - \theta_{0,K(r)}) \\ &\quad + \frac{1}{3!} \sum_{1 \leq i, j, l \leq K-r} \left( \frac{\partial^3 \ell(\theta^\Delta)}{\partial \theta_i \partial \theta_j \partial \theta_l} \right) (\hat{\theta}_{K(r),i} - \theta_{0,i}) (\hat{\theta}_{K(r),j} - \theta_{0,j}) (\hat{\theta}_{K(r),l} - \theta_{0,l}), \quad (12)\end{aligned}$$

where  $\|\theta^\Delta - \theta_0^*\| \leq \|\hat{\theta}^* - \theta_0^*\|$  and  $\frac{\partial \ell(\theta)}{\partial \theta_{K(r)}} = \left( \frac{\partial \ell(\theta)}{\partial \theta_1}, \dots, \frac{\partial \ell(\theta)}{\partial \theta_{K-r}} \right)'$ , and in general

$$\begin{aligned}\ell(\hat{\theta}) &= \ell(\theta_0) + \left( \frac{\partial \ell(\theta_0)}{\partial \theta} \right)' (\hat{\theta} - \theta_0) \\ &\quad + \frac{1}{2} (\hat{\theta} - \theta_0)' \left( \frac{\partial^2 \ell(\theta_0)}{\partial \theta_i \partial \theta_j} \right) (\hat{\theta} - \theta_0) \\ &\quad + \frac{1}{3!} \sum_{1 \leq i, j, l \leq K} \left( \frac{\partial^3 \ell(\theta^\nabla)}{\partial \theta_i \partial \theta_j \partial \theta_l} \right) (\hat{\theta}_i - \theta_{0,i}) (\hat{\theta}_j - \theta_{0,j}) (\hat{\theta}_l - \theta_{0,l}), \quad (13)\end{aligned}$$

where  $\|\theta^\nabla - \theta_0\| \leq \|\hat{\theta} - \theta_0\|$ .

- By the CLT for  $\hat{\boldsymbol{\theta}}_{K(r)}$  and  $\hat{\boldsymbol{\theta}}$  and assuming there exists small  $\delta > 0$  such that

$$\sup_{\boldsymbol{\theta} \in B_\delta(\boldsymbol{\theta}_0^*)} \left| \frac{1}{n} \frac{\partial^3 \ell(\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j \partial \theta_l} \right| = O_p(1) \text{ for } 1 \leq i, j, l \leq K,$$

it can be shown that under  $H_0$ , the third terms on the RHS of (12) and (13) are of order  $o_p(1)$ .

- Therefore, it remains to show that

$$(I) - (II) \xrightarrow[\text{H}_0]{d} \chi^2(r), \quad (14)$$

where

$$\begin{aligned} (I) &= 2 \left( \frac{\partial \ell(\boldsymbol{\theta}_0^*)}{\partial \boldsymbol{\theta}} \right)' (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0^*) + (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0^*)' \left( \frac{\partial^2 \ell(\boldsymbol{\theta}_0^*)}{\partial \theta_i \partial \theta_j} \right) (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0^*), \\ (II) &= 2 \left( \frac{\partial \ell(\boldsymbol{\theta}_0^*)}{\partial \boldsymbol{\theta}_{K(r)}} \right)' (\hat{\boldsymbol{\theta}}_{K(r)} - \boldsymbol{\theta}_{0,K(r)}) \\ &\quad + (\hat{\boldsymbol{\theta}}_{K(r)} - \boldsymbol{\theta}_{0,K(r)})' \left( \frac{\partial^2 \ell(\boldsymbol{\theta}_0^*)}{\partial \theta_i \partial \theta_j} \right)_{1 \leq i, j \leq K-r} (\hat{\boldsymbol{\theta}}_{K(r)} - \boldsymbol{\theta}_{0,K(r)}). \end{aligned}$$

Define

$$\mathbf{x}_t = \frac{\partial \log f_{\theta_0^*}(x_t)}{\partial \theta} \quad \text{and} \quad \mathbf{x}_{t,K(r)} = \frac{\partial \log f_{\theta_0^*}(x_t)}{\partial \theta_{K(r)}}.$$

Then, under the regularity conditions similar to those used to obtain the CLTs for  $\hat{\theta}$  and  $\hat{\theta}_{K(r)}$ , we have

$$\begin{aligned} \text{(I)} &= \mathbf{1}' \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{1} + o_p(1), \\ \text{(II)} &= \mathbf{1}' \mathbf{X}_{K(r)} (\mathbf{X}'_{K(r)} \mathbf{X}_{K(r)})^{-1} \mathbf{X}'_{K(r)} \mathbf{1} + o_p(1), \end{aligned} \tag{15}$$

where

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}'_1 \\ \vdots \\ \mathbf{x}'_n \end{pmatrix}, \mathbf{X}_{K(r)} = \begin{pmatrix} \mathbf{x}'_{1,K(r)} \\ \vdots \\ \mathbf{x}'_{n,K(r)} \end{pmatrix}, \text{ and } \mathbf{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}.$$

Let  $\mathbf{X}_{-1}$  satisfy  $(\mathbf{X}_{K(r)}, \mathbf{X}_{-1}) = \mathbf{X}$  and define

$$\mathbf{D} = \begin{pmatrix} \mathbf{I}_{(K-r) \times (K-r)} & -(\mathbf{X}'_{K(r)} \mathbf{X}_{K(r)})^{-1} \mathbf{X}'_{K(r)} \mathbf{X}_{-1} \\ \mathbf{0}_{r \times (K-r)} & \mathbf{I}_{r \times r} \end{pmatrix}.$$

Then

$$\begin{aligned} & \mathbf{1}' \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{1} \\ = & \mathbf{1}' \mathbf{X} \mathbf{D} (\mathbf{D}' \mathbf{X}' \mathbf{X} \mathbf{D})^{-1} \mathbf{D}' \mathbf{X}' \mathbf{1} \\ \stackrel{\text{why?}}{=} & \mathbf{1}' \mathbf{X}_{K(r)} (\mathbf{X}'_{K(r)} \mathbf{X}_{K(r)})^{-1} \mathbf{X}'_{K(r)} \mathbf{1} \\ & + \mathbf{1}' (\mathbf{I} - \mathbf{M}_{K(r)}) \mathbf{X}_{-1} (\mathbf{X}'_{-1} (\mathbf{I} - \mathbf{M}_{K(r)}) \mathbf{X}_{-1})^{-1} \mathbf{X}'_{-1} (\mathbf{I} - \mathbf{M}_{K(r)}) \mathbf{1}, \quad (16) \end{aligned}$$

where

$$\mathbf{M}_{K(r)} = \mathbf{X}_{K(r)} (\mathbf{X}'_{K(r)} \mathbf{X}_{K(r)})^{-1} \mathbf{X}'_{K(r)}$$

is the orthogonal projection matrix onto  $C(\mathbf{X}_{K(r)})$ .

Since

$$\frac{1}{n} \mathbf{X}' \mathbf{X} \xrightarrow[\text{H}_0]{pr.} \mathbf{I}(\boldsymbol{\theta}_0^*) = \begin{pmatrix} \mathbf{I}_{11}(\boldsymbol{\theta}_0^*) & \mathbf{I}_{12}(\boldsymbol{\theta}_0^*) \\ \mathbf{I}_{21}(\boldsymbol{\theta}_0^*) & \mathbf{I}_{22}(\boldsymbol{\theta}_0^*) \end{pmatrix} \quad (\text{p.d.})$$

and

$$\frac{1}{\sqrt{n}} \sum_{t=1}^n \mathbf{x}_t \xrightarrow[\text{H}_0]{d} N(\mathbf{0}, \mathbf{I}(\boldsymbol{\theta}_0^*))$$

under some regularity conditions, it holds that

$$\begin{aligned} \frac{1}{n} \mathbf{X}'_{-1} (\mathbf{I} - \mathbf{M}_{K(r)}) \mathbf{X}_{-1} &\xrightarrow[\text{H}_0]{pr.} \mathbf{I}_{22}(\boldsymbol{\theta}_0^*) - \mathbf{I}_{21}(\boldsymbol{\theta}_0^*) \mathbf{I}_{11}^{-1}(\boldsymbol{\theta}_0^*) \mathbf{I}_{12}(\boldsymbol{\theta}_0^*), \\ \frac{1}{\sqrt{n}} \mathbf{X}'_{-1} (\mathbf{I} - \mathbf{M}_{K(r)}) \mathbf{1} &\stackrel{\text{why?}}{=} (-\mathbf{X}'_{-1} \mathbf{X}_{K(r)} (\mathbf{X}'_{K(r)} \mathbf{X}_{K(r)})^{-1}, \mathbf{I}_{r \times r}) \frac{1}{\sqrt{n}} \sum_{t=1}^n \mathbf{x}_t \\ &\stackrel{\text{H}_0}{=} (-\mathbf{I}_{21}(\boldsymbol{\theta}_0^*) \mathbf{I}_{11}^{-1}(\boldsymbol{\theta}_0^*), \mathbf{I}_{r \times r}) \frac{1}{\sqrt{n}} \sum_{t=1}^n \mathbf{x}_t + o_p(1), \end{aligned}$$

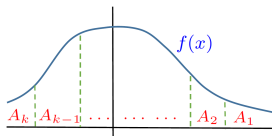
and

$$\mathbf{1}' (\mathbf{I} - \mathbf{M}_{K(r)}) \mathbf{X}_{-1} (\mathbf{X}'_{-1} (\mathbf{I} - \mathbf{M}_{K(r)}) \mathbf{X}_{-1})^{-1} \mathbf{X}'_{-1} (\mathbf{I} - \mathbf{M}_{K(r)}) \mathbf{1} \xrightarrow[\text{H}_0]{d} \chi^2(r). \quad (17)$$

Now, the desired conclusion (14) follows from (15)–(17).

# Pearson's Chi-Squared Test with unknown parameters

- Assume  $x_1, \dots, x_n$  is a random sample generated from  $f(x)$ . Consider



- Let  $p_i = P(x_1 \in A_i) = \int_{z \in A_i} f(z) dz$ , where  $\bigcup_{i=1}^k A_i = \mathbb{R}$  and  $A_i \cap A_j = \emptyset$  if  $i \neq j$ .
- Define

$$\mathbf{y}_i = \begin{pmatrix} I_{\{x_i \in A_1\}} \\ \vdots \\ I_{\{x_i \in A_{k-1}\}} \end{pmatrix} \equiv \begin{pmatrix} y_{i,1} \\ \vdots \\ y_{i,k-1} \end{pmatrix}.$$

Then, the likelihood function of  $\mathbf{y}_1, \dots, \mathbf{y}_n$  is

$$L(p_1, \dots, p_{k-1}) = \prod_{i=1}^n p_1^{y_{i,1}} \cdots p_k^{y_{i,k}} = p_1^{\sum_{i=1}^n y_{i,1}} \cdots p_k^{\sum_{i=1}^n y_{i,k}} \equiv p_1^{O_1} \cdots p_k^{O_k},$$

where  $y_{i,k} = 1 - \sum_{j=1}^{k-1} y_{i,j}$  and  $p_k = 1 - \sum_{j=1}^{k-1} p_j$ .

- Now Consider the null hypothesis,

$$H_0 : p_i = \int_{z \in A_i} f_{\theta}(z) dz = p_i(\theta), \quad i = 1, \dots, k-1,$$

versus the alternative,

$$H_A : H_0 \text{ is wrong,}$$

where  $f_{\theta}(\cdot)$  is a probability density function indexed by  $\theta \in \Theta$ .

- Then, one can use the likelihood ratio statistics

$$-2 \log \Lambda_n \xrightarrow{d} \chi^2(k-1 - \dim(\Theta)),$$

where

$$\Lambda_n = \frac{\sup_{\theta \in \Theta} L(p_1(\theta), \dots, p_{k-1}(\theta))}{\sup_{\sum_{i=1}^k p_i=1} L(p_1, \dots, p_{k-1})} = \frac{\prod_{j=1}^k (p_j(\hat{\theta}))^{O_j}}{\prod_{j=1}^k \hat{p}_j^{O_j}},$$

where  $\hat{p}_j = O_j/n$ , and reject  $H_0$  if  $-2 \log \Lambda_n$  is large.

- In the following, I'll show that

$$\sum_{j=1}^k \frac{(O_j - np_j(\hat{\theta}))^2}{np_j(\hat{\theta})} \xrightarrow[H_0]{d} \chi^2(k-1 - \dim(\Theta)).$$



Note that under  $H_0$ ,

$$\begin{aligned}
 -2 \log \Lambda_n &= 2 \sum_{j=1}^k O_j \log \left( \frac{\hat{p}_j}{p_j(\hat{\theta})} - 1 + 1 \right) \\
 &\sim 2 \sum_{j=1}^k O_j \left\{ \frac{\hat{p}_j - p_j(\hat{\theta})}{p_j(\hat{\theta})} - \frac{1}{2} \left[ \frac{\hat{p}_j - p_j(\hat{\theta})}{p_j(\hat{\theta})} \right]^2 \right\} \\
 &= 2 \sum_{j=1}^k O_j \left\{ \frac{\hat{p}_j - p_j(\hat{\theta})}{\hat{p}_j} \left( \frac{\hat{p}_j}{p_j(\hat{\theta})} - 1 + 1 \right) - \frac{1}{2} \left[ \frac{(\hat{p}_j - p_j(\hat{\theta}))^2}{\hat{p}_j p_j(\hat{\theta})} \left( \frac{\hat{p}_j}{p_j(\hat{\theta})} - 1 + 1 \right) \right] \right\} \\
 &= 2 \sum_{j=1}^k O_j \frac{\hat{p}_j - p_j(\hat{\theta})}{\hat{p}_j} + \sum_{j=1}^k O_j \frac{(\hat{p}_j - p_j(\hat{\theta}))^2}{\hat{p}_j p_j(\hat{\theta})} - \sum_{j=1}^k O_j \frac{(\hat{p}_j - p_j(\hat{\theta}))^2}{\hat{p}_j p_j(\hat{\theta})} \left( \frac{\hat{p}_j - p_j(\hat{\theta})}{p_j(\hat{\theta})} \right) \\
 &= 2n \sum_{j=1}^k (\hat{p}_j - p_j(\hat{\theta})) + \sum_{j=1}^k O_j \frac{(\hat{p}_j - p_j(\hat{\theta}))^2}{\hat{p}_j p_j(\hat{\theta})} - \sum_{j=1}^k O_j \frac{(\hat{p}_j - p_j(\hat{\theta}))^2}{\hat{p}_j p_j(\hat{\theta})} \left( \frac{\hat{p}_j - p_j(\hat{\theta})}{p_j(\hat{\theta})} \right) \\
 &= \sum_{j=1}^k O_j \frac{(\hat{p}_j - p_j(\hat{\theta}))^2}{\hat{p}_j p_j(\hat{\theta})} - \sum_{j=1}^k O_j \frac{(\hat{p}_j - p_j(\hat{\theta}))^2}{\hat{p}_j p_j(\hat{\theta})} \left( \frac{\hat{p}_j - p_j(\hat{\theta})}{p_j(\hat{\theta})} \right) \\
 &= \sum_{j=1}^k O_j \frac{(\hat{p}_j - p_j(\hat{\theta}))^2}{\hat{p}_j p_j(\hat{\theta})} (1 + o_p(1)) \quad \left( \text{under } H_0, \max_{1 \leq i \leq k} \left| \frac{\hat{p}_i - p_i(\hat{\theta})}{p_i(\hat{\theta})} \right| = o_p(1) \right) \\
 &= \sum_{j=1}^k \frac{(O_j - np_j(\hat{\theta}))^2}{np_j(\hat{\theta})} (1 + o_p(1)).
 \end{aligned}$$

- Hence, by the following Lemma,

$$-2 \log \Lambda_n - \sum_{j=1}^k \frac{(O_j - np_j(\hat{\boldsymbol{\theta}}))^2}{np_j(\hat{\boldsymbol{\theta}})} = o_p(1).$$

**Lemma:** If  $A_n = B_n(1 + o_p(1))$  and  $A_n = O_p(1)$ , then  $A_n - B_n \xrightarrow{pr.} 0$ .

**Proof.** Note that  $A_n - B_n = (B_n - A_n + A_n)o_p(1) = (B_n - A_n)o_p(1) + o_p(1)$ , which implies  $(A_n - B_n)(1 + o_p(1)) = o_p(1)$ , and hence  $A_n - B_n = \frac{o_p(1)}{1 + o_p(1)} = o_p(1)$ .

- By Slutsky's Theorem,

$$\sum_{j=1}^k \frac{(O_j - np_j(\hat{\boldsymbol{\theta}}))^2}{np_j(\hat{\boldsymbol{\theta}})} \xrightarrow[H_0]{d} \chi^2(k - 1 - \dim(\boldsymbol{\Theta})).$$

Therefore,  $-2 \log \Lambda_n$  and  $\sum_{j=1}^k \frac{(O_j - np_j(\hat{\boldsymbol{\theta}}))^2}{np_j(\hat{\boldsymbol{\theta}})}$  have the same limiting distribution.

## Application

Assume that there are  $r$  rows and  $s$  columns in the contingency table, and the total number in the table is  $n$ . Let  $O_{ij}$  be the observed frequency for the  $i^{\text{th}}$  row and  $j^{\text{th}}$  column, and  $p_{ij}$  be the probability that an observation falls into the  $i^{\text{th}}$  row and  $j^{\text{th}}$  column. Define  $p_{i\cdot} = \sum_{j=1}^s p_{ij}$  and  $p_{\cdot j} = \sum_{i=1}^r p_{ij}$ . For the hypothesis

$$H_0 : p_{ij} = p_{i\cdot} p_{\cdot j} \text{ for } i = 1, \dots, r \text{ and } j = 1, \dots, s, \quad \text{versus} \quad H_A : H_0 \text{ is wrong,}$$

the test statistic is

$$\sum_{i=1}^r \sum_{j=1}^s \frac{(O_{ij} - n\hat{p}_{i\cdot}\hat{p}_{\cdot j})^2}{n\hat{p}_{i\cdot}\hat{p}_{\cdot j}} \xrightarrow[H_0]{d} \chi^2((r-1)(s-1)),$$

where  $\hat{p}_{i\cdot} = \sum_{j=1}^s O_{ij}/n$  and  $\hat{p}_{\cdot j} = \sum_{i=1}^r O_{ij}/n$ .

## Remark

$$rs - 1 - (r - 1) - (s - 1) = rs - 1 - r + 1 - s + 1 = rs - r - s + 1 = (r - 1)(s - 1).$$

# Akaike's information criterion (AIC)

- Let

$g(\mathbf{y}^*) = g(y_1^*, \dots, y_n^*)$  denote the likelihood function of  $\mathbf{y}^*$  (true pdf of  $\mathbf{y}^*$ )

and

$f_{\theta}(\mathbf{y}^*)$  is a family of approximation models indexed by  $\theta$  (a family of pdfs indexed by  $\theta$ ).

- The "distance" between  $g(\mathbf{y}^*)$  and  $f_{\theta}(\mathbf{y}^*)$  is measured by Kullback-Leibler distance

$$\text{KL}(f_{\theta}, g) = -E_g \left( \log \frac{f_{\theta}(\mathbf{y}^*)}{g(\mathbf{y}^*)} \right) = -E_g(\log f_{\theta}(\mathbf{y}^*)) + E_g(\log g(\mathbf{y}^*)). \quad (18)$$

## Remark

$\text{KL}(g, g) = 0$  and  $\text{KL}(f_{\theta}, g) \geq 0$ , since  $-\log t$  is a convex function, leading to

$$E(-\log X) \geq -\log E(X)$$

by Jensen's inequality.

- More specifically, letting

$$X = \frac{f_{\theta}(\mathbf{y}^*)}{g(\mathbf{y}^*)},$$

we have by Jensen's inequality and the convexity of  $-\log t$ ,

$$E_g \left( -\log \frac{f_{\theta}(\mathbf{y}^*)}{g(\mathbf{y}^*)} \right) \geq -\log E_g \left( \frac{f_{\theta}(\mathbf{y}^*)}{g(\mathbf{y}^*)} \right).$$

- Moreover,

$$E_g \left( \frac{f_{\theta}(\mathbf{y}^*)}{g(\mathbf{y}^*)} \right) = \int \frac{f_{\theta}(\mathbf{y}^*)}{g(\mathbf{y}^*)} g(\mathbf{y}^*) d\mathbf{y}^* = \int f_{\theta}(\mathbf{y}^*) d\mathbf{y}^* = 1$$

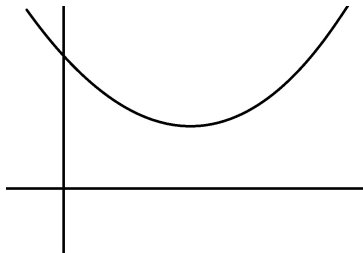
because  $f_{\theta}(\mathbf{y}^*)$  is a pdf and hence integrates to "1".

Note that Jensen's inequality says that for any convex function  $h(\cdot)$ , and random variable  $X$ ,

$$E(h(X)) \geq h(E(X)),$$

whenever the expectations exist. For example,

$$\begin{aligned} E(X^2) &\geq (EX)^2, \\ -E(\log X) &\geq -\log EX, \\ E(e^X) &\geq e^{E(X)}, \\ &\vdots \end{aligned}$$



- In the derivation of AIC,  $\mathbf{y}^*$  is regarded as the "future", whereas  $\mathbf{y}$ , an independent copy of  $\mathbf{y}^*$  (namely  $\mathbf{y}$  and  $\mathbf{y}^*$  are independent, but have the same distribution.), denotes the observations at hand.
- Since the factor  $E_g(\log g(\mathbf{y}^*))$  in (18) will be cancelled out when comparing across different approximation models (families), the goal here is to construct an asymptotically unbiased estimate of

$$-E_g(\log f_{\hat{\theta}(\mathbf{y})}(\mathbf{y}^*)), \quad (g: \text{true model}; \hat{\theta}(\mathbf{y}): \text{present}; \mathbf{y}^*: \text{future})$$

in which  $\hat{\theta}(\mathbf{y})$  is the MLE of  $\theta$  using observations  $\mathbf{y}$ .

- A natural estimate of  $-E_g(\log f_{\hat{\theta}(\mathbf{y})}(\mathbf{y}^*))$  is given by

$$-\log f_{\hat{\theta}(\mathbf{y})}(\mathbf{y}), \quad (g \rightarrow \text{unknown} \rightarrow \text{dropped here}; \mathbf{y}^* \text{ is replaced by } \mathbf{y})$$

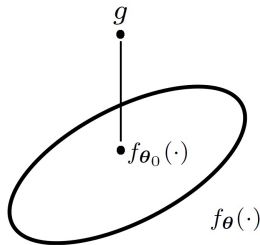
the log likelihood function of  $\mathbf{y}$  based on  $f_{\theta}(\cdot)$ .

- However, a "bias correction" term is needed to achieve asymptotic unbiasedness.

To find the bias correction term, note first that by Taylor's theorem,

$$\begin{aligned} E_g(\log f_{\hat{\theta}(\mathbf{y})}(\mathbf{y}^*)) &= E_g(\log f_{\theta_0}(\mathbf{y}^*)) \\ &+ E_g \left[ \left( \frac{\partial \log f_{\theta_0}(\mathbf{y}^*)}{\partial \theta} \right)' (\hat{\theta}(\mathbf{y}) - \theta_0) \right] \\ &+ \frac{1}{2} E_g \left[ (\hat{\theta}(\mathbf{y}) - \theta_0)' \left( \frac{\partial^2 \log f_{\theta_0}(\mathbf{y}^*)}{\partial \theta_i \partial \theta_j} \right) (\hat{\theta}(\mathbf{y}) - \theta_0) \right], \quad (19) \end{aligned}$$

where  $\theta_0 = \underset{\theta \in \Theta}{\operatorname{argmin}} \operatorname{KL}(f_{\theta}, g)$  and  $\|\theta^* - \theta_0\| \leq \|\hat{\theta}(\mathbf{y}) - \theta_0\|$ .





Assume  $f_{\theta_0} = g$  (the true model is included among the approximation family). (19) becomes

$$\begin{aligned}
 & E_g(\log f_{\hat{\theta}(\mathbf{y})}(\mathbf{y}^*)) \\
 = & E_{f_{\theta_0}}(\log f_{\theta_0}(\mathbf{y}^*)) + E_{f_{\theta_0}} \left[ \left( \frac{\partial \log f_{\theta_0}(\mathbf{y}^*)}{\partial \boldsymbol{\theta}} \right)' (\hat{\boldsymbol{\theta}}(\mathbf{y}) - \boldsymbol{\theta}_0) \right] \\
 & + \frac{1}{2} E_{f_{\theta_0}} \left[ (\hat{\boldsymbol{\theta}}(\mathbf{y}) - \boldsymbol{\theta}_0)' \frac{\partial^2 \log f_{\theta_0}(\mathbf{y}^*)}{\partial \theta_i \partial \theta_j} (\hat{\boldsymbol{\theta}}(\mathbf{y}) - \boldsymbol{\theta}_0) \right] \\
 \sim & E_{f_{\theta_0}}(\log f_{\theta_0}(\mathbf{y}^*)) + \frac{1}{2} E_{f_{\theta_0}} \left[ (\hat{\boldsymbol{\theta}}(\mathbf{y}) - \boldsymbol{\theta}_0)' \frac{\partial^2 \log f_{\theta_0}(\mathbf{y}^*)}{\partial \theta_i \partial \theta_j} (\hat{\boldsymbol{\theta}}(\mathbf{y}) - \boldsymbol{\theta}_0) \right] \quad (20)
 \end{aligned}$$

because  $\boldsymbol{\theta}^* \sim \boldsymbol{\theta}_0$  and

$$\begin{aligned}
 & E_{f_{\theta_0}} \left[ \left( \frac{\partial \log f_{\theta_0}(\mathbf{y}^*)}{\partial \boldsymbol{\theta}} \right)' (\hat{\boldsymbol{\theta}}(\mathbf{y}) - \boldsymbol{\theta}_0) \right] \\
 = & E_{f_{\theta_0}} \left( \frac{\partial \log f_{\theta_0}(\mathbf{y}^*)}{\partial \boldsymbol{\theta}} \right)' E(\hat{\boldsymbol{\theta}}(\mathbf{y}) - \boldsymbol{\theta}_0) \quad (\mathbf{y}^* \text{ and } \mathbf{y} \text{ are independent}) \\
 = & 0. \quad \left( E_{f_{\theta_0}} \left( \frac{\partial \log f_{\theta_0}(\mathbf{y}^*)}{\partial \boldsymbol{\theta}} \right) = \mathbf{0} \right)
 \end{aligned}$$

- Moreover, assume

$$\lim_{n \rightarrow \infty} \frac{1}{n} E_{f_{\theta_0}} \left( -\frac{\partial^2 \log f_{\theta_0}(\mathbf{y}^*)}{\partial \theta_i \partial \theta_j} \right) = \mathbf{I}(\theta_0) \text{ exist. } (\mathbf{y}^* \text{ can also be replaced by } \mathbf{y})$$

- Then, we have shown (heuristically) that

$$\sqrt{n}(\hat{\theta}(\mathbf{y}) - \theta_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{I}^{-1}(\theta_0)),$$

and

$$-\frac{1}{n} \frac{\partial^2 \log f_{\theta_0}(\mathbf{y}^*)}{\partial \theta_i \partial \theta_j} \xrightarrow{pr.} \mathbf{I}(\theta_0),$$

leading to

$$-(\hat{\theta}(\mathbf{y}) - \theta_0)' \frac{\partial^2 \log f_{\theta_0}(\mathbf{y}^*)}{\partial \theta_i \partial \theta_j} (\hat{\theta}(\mathbf{y}) - \theta_0) \xrightarrow{d} \chi^2(p),$$

assuming  $\theta$  is a  $p$ -dimensional vector.

- As a result, (20) becomes

$$\begin{aligned} & E_{f_{\theta_0}}(\log f_{\hat{\theta}(\mathbf{y})}(\mathbf{y}^*)) \\ \sim & E_{f_{\theta_0}}(\log f_{\theta_0}(\mathbf{y}^*)) - \frac{p}{2} \\ = & E_{f_{\theta_0}}(\log f_{\theta_0}(\mathbf{y})) - \frac{p}{2} \quad (\text{this replacement won't change anything}) \end{aligned} \quad (21)$$

- By Taylor's expansion again,

$$\begin{aligned}\log f_{\theta_0}(\mathbf{y}) &= \log f_{\hat{\theta}(\mathbf{y})}(\mathbf{y}) + \left( \frac{\partial \log f_{\hat{\theta}(\mathbf{y})}(\mathbf{y})}{\partial \theta} \right)' (\theta_0 - \hat{\theta}(\mathbf{y})) \\ &\quad + \frac{1}{2} (\hat{\theta}(\mathbf{y}) - \theta_0)' \left( \frac{\partial^2 \log f_{\hat{\theta}(\mathbf{y})}(\mathbf{y})}{\partial \theta_i \partial \theta_j} \right) (\hat{\theta}(\mathbf{y}) - \theta_0),\end{aligned}$$

where  $\|\hat{\theta} - \theta_0\| \leq \|\hat{\theta}(\mathbf{y}) - \theta_0\|$ .

- Since

$$\frac{\partial \log f_{\hat{\theta}(\mathbf{y})}(\mathbf{y})}{\partial \theta} = \mathbf{0} \quad (\text{why?})$$

and

$$\begin{aligned}& -(\hat{\theta}(\mathbf{y}) - \theta_0)' \frac{\partial^2 \log f_{\hat{\theta}(\mathbf{y})}(\mathbf{y})}{\partial \theta_i \partial \theta_j} (\hat{\theta}(\mathbf{y}) - \theta_0) \\ \sim & -(\hat{\theta}(\mathbf{y}) - \theta_0)' \frac{\partial^2 \log f_{\theta_0}(\mathbf{y})}{\partial \theta_i \partial \theta_j} (\hat{\theta}(\mathbf{y}) - \theta_0) \xrightarrow{d} \chi^2(p),\end{aligned}$$

it follows that

$$E_{f_{\theta_0}}(\log f_{\theta_0}(\mathbf{y})) \sim E_{f_{\theta_0}}(\log f_{\hat{\theta}(\mathbf{y})}(\mathbf{y})) - \frac{p}{2}. \quad (22)$$

- Combining (21) and (22) yields

$$-E_{f_{\theta_0}}(\log f_{\hat{\theta}(\mathbf{y})}(\mathbf{y}^*)) \sim E_{f_{\theta_0}}(-\log f_{\hat{\theta}(\mathbf{y})}(\mathbf{y}) + p). \quad (23)$$

This is what we call "asymptotic unbiasedness". ( $p$ : This is our bias correction term!!)

- Now, the definition of AIC is given by

$$-2 \log f_{\hat{\theta}(\mathbf{y})} + 2p,$$

which is the quantity inside the expectation on the RHS of (23) multiplied by "2".

# More on information criteria

Let

$$f_{\theta}(\mathbf{y}, \mathbf{x}_1, \dots, \mathbf{x}_n) = (2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}'_i \boldsymbol{\beta})^2 \right\}, \quad (24)$$

where  $\boldsymbol{\theta} = (\boldsymbol{\beta}', \sigma^2)'$ . Then, the MLE of  $\boldsymbol{\theta}$  is  $\hat{\boldsymbol{\theta}} = (\hat{\boldsymbol{\beta}}', \hat{\sigma}^2)'$ , where

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}'_i \hat{\boldsymbol{\beta}})^2, \quad \text{and} \quad \mathbf{X} = \begin{pmatrix} \mathbf{x}'_1 \\ \vdots \\ \mathbf{x}'_n \end{pmatrix},$$

yielding

$$\begin{aligned} \log f_{\hat{\boldsymbol{\theta}}}(\cdot) &= -\frac{n}{2}(\log(2\pi) + \log \hat{\sigma}^2) - \frac{1}{2\hat{\sigma}^2} \sum_{i=1}^n (y_i - \mathbf{x}'_i \hat{\boldsymbol{\beta}})^2 \\ &= -\frac{n}{2} \log(2\pi) - \frac{n}{2} - \frac{n}{2} \log \hat{\sigma}^2. \end{aligned}$$

( $-\frac{n}{2} \log(2\pi) - \frac{n}{2}$ : this part is independent of the model)

- Therefore, the essential part of AIC is

$$n \log \hat{\sigma}^2 + 2\dim(\beta), \quad (25)$$

noting that the number of parameters in model (24) is  $\dim(\beta) + 1$  in which "1" will be cancelled out when model comparisons are performed.

- In most practical situations, the distribution of error in model (24) is unknown. Therefore, AIC is in general not obtainable. But we can still "borrow" the AIC obtained in the Gaussian case, namely (25), to do model selection.

# Bayesian information criterion (BIC)

- Bayesian information criterion (BIC):

$$-2 \log f_{\hat{\theta}}(\mathbf{y}) + \log n \#(\text{estimated parameters})$$

- $\log n$ : in contrast to "2" for AIC
- $\#(\text{estimated parameters})$ : the number of meaning estimated parameters
- BIC in linear regression model with Gaussian error:

$$n \log \hat{\sigma}^2 + \log n \dim(\beta).$$

# Consistency of BIC in regression models

- Consider a linear model

$$y_t = \mathbf{x}_t' \boldsymbol{\beta} + \epsilon_t, \quad \epsilon_t \stackrel{\text{i.i.d.}}{\sim} (0, \sigma^2)$$

where  $\mathbf{x}_t = (x_{t1}, \dots, x_{tK})'$  and  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_K)'$  in which some of  $\beta_i$  may equal 0.

- Define  $J_0 = \{\beta_i : 1 \leq i \leq K, \beta_i \neq 0\}$  and

$$\text{BIC}(J) = n \log \hat{\sigma}_J^2 + \log n \times \#(J),$$

where  $J$  is a subset of  $\mathcal{K} = \{1, \dots, K\}$ ,  $\#(J)$  is **the number of elements in  $J$** ,

$$\hat{\sigma}_J^2 = \frac{1}{n} \sum_{t=1}^n (y_t - \mathbf{x}_t'(J) \hat{\boldsymbol{\beta}}(J))^2$$

is **the residual mean squared error of model  $J$** , with  $\mathbf{x}_t(J) = (x_{ti}, i \in J)$  and

$$\hat{\boldsymbol{\beta}}(J) = \left( \sum_{t=1}^n \mathbf{x}_t(J) \mathbf{x}_t'(J) \right)^{-1} \sum_{t=1}^n \mathbf{x}_t(J) y_t$$

denoting the LSE of model  $J$ .



- Define

$$\hat{J} = \operatorname{argmin}_{J \subseteq 2^{\mathcal{K}}} \text{BIC}(J),$$

where  $2^{\mathcal{K}}$  denotes all subsets of  $\mathcal{K}$ .

- In the following, I shall show that

$$\lim_{n \rightarrow \infty} P(\hat{J} = J_0) = 1, \quad (J_0 : \text{true model}) \quad (26)$$

"without" assuming that  $\epsilon_t$  are Gaussian.

# Proof

- If we can show that for  $J_0 - J = \emptyset$  and  $\#(J) > \#(J_0)$  ( $J$  嚴格包含  $J_0$ ),

$$P(\text{BIC}(J) \leq \text{BIC}(J_0)) \xrightarrow{n \rightarrow \infty} 0, \quad (27)$$

and for  $J_0 - J \neq \emptyset$  ( $J$  不包含  $J_0$ ),

$$P(\text{BIC}(J) \leq \text{BIC}(J_0)) \xrightarrow{n \rightarrow \infty} 0, \quad (28)$$

then combining (27) and (28) leads to the desired conclusion (26).

- To show (27), note first that

$$\{\text{BIC}(J) \leq \text{BIC}(J_0)\} = \{n(\log \hat{\sigma}_{J_0}^2 - \log \hat{\sigma}_J^2) \geq \log n(\#(J) - \#(J_0))\}. \quad (29)$$

Moreover, since

$$\begin{aligned}
 n(\log \hat{\sigma}_{J_0}^2 - \log \hat{\sigma}_J^2) &= n \left( \log \left( 1 + \frac{\hat{\sigma}_{J_0}^2 - \hat{\sigma}_J^2}{\hat{\sigma}_J^2} \right) \right) \sim \frac{n}{\sigma^2} (\hat{\sigma}_{J_0}^2 - \hat{\sigma}_J^2) \\
 &= \frac{1}{\sigma^2} \epsilon' (M_J - M_{J_0}) \epsilon \\
 &\leq \frac{1}{\sigma^2} \epsilon' M_J \epsilon,
 \end{aligned}$$

and since

$$\frac{1}{\sigma^2} E(\epsilon' M_J \epsilon) = \#(J) < \infty,$$

we have

$$n(\log \hat{\sigma}_{J_0}^2 - \log \hat{\sigma}_J^2) = O_p(1). \quad (\text{why?}) \quad (30)$$

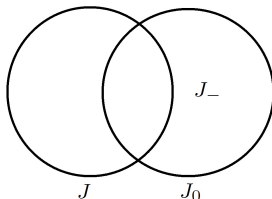
- $\sim$  holds due to "correctness" of those two models ( $\hat{\sigma}_{J_0}^2 \xrightarrow{pr.} \sigma^2$  and  $\hat{\sigma}_J^2 \xrightarrow{pr.} \sigma^2$ ) and by Taylor's expansion.
- $M_J$  and  $M_{J_0}$  are orthogonal projection matrices for the column spaces of

$$\mathbf{X}_J = \begin{pmatrix} \mathbf{x}'_1(J) \\ \vdots \\ \mathbf{x}'_n(J) \end{pmatrix} \quad \text{and} \quad \mathbf{X}_{J_0} = \begin{pmatrix} \mathbf{x}'_1(J_0) \\ \vdots \\ \mathbf{x}'_n(J_0) \end{pmatrix}.$$

- Now, (27) follows directly from (29), (30) and  $\log n \rightarrow \infty$ .
- To show (28), define

$$J^* = J \cup J_0, \text{ and } J_- = J_0 - J,$$

noting that  $J^* = J \cup J_-$  and  $J \cap J_- = \emptyset$ .



- In view of (27), (28) is guaranteed by

$$P(\text{BIC}(J) \leq \text{BIC}(J^*)) \xrightarrow{n \rightarrow \infty} 0. \quad (31)$$

- It is clear that

$$\begin{aligned} \{\text{BIC}(J) \leq \text{BIC}(J^*)\} &= \{n(\log \hat{\sigma}_J^2 - \log \hat{\sigma}_{J^*}^2) \leq \log n(\#(J^*) - \#(J))\}, \\ \text{and } \log \hat{\sigma}_{J^*}^2 &\xrightarrow{pr.} \log \sigma^2. \end{aligned} \quad (32)$$

(Since  $J^*$  is a correct model, namely, a model including  $J_0$  as a subset model)

In addition, we have

$$\mathbf{y} = \mathbf{X}_{J^*} \boldsymbol{\beta}_{J^*} + \boldsymbol{\epsilon}, \quad (\text{why?})$$

where  $\boldsymbol{\beta}_{J^*} = (\beta_j, j \in J^*)$ , and hence

$$\begin{aligned} \hat{\sigma}_J^2 &= \frac{1}{n} \mathbf{y}' (\mathbf{I} - \mathbf{M}_J) \mathbf{y} \\ &= \frac{1}{n} \mathbf{y}' (\mathbf{I} - \mathbf{M}_J) (\mathbf{I} - \mathbf{M}_J) (\mathbf{X}_J \boldsymbol{\beta}_J + \mathbf{X}_{J_-} \boldsymbol{\beta}_{J_-} + \boldsymbol{\epsilon}) \\ &= \frac{1}{n} \mathbf{y}' (\mathbf{I} - \mathbf{M}_J) (\mathbf{I} - \mathbf{M}_J) (\mathbf{X}_{J_-} \boldsymbol{\beta}_{J_-} + \boldsymbol{\epsilon}) \\ &= \frac{1}{n} (\mathbf{X}_{J_-} \boldsymbol{\beta}_{J_-} + \boldsymbol{\epsilon})' (\mathbf{I} - \mathbf{M}_J) (\mathbf{X}_{J_-} \boldsymbol{\beta}_{J_-} + \boldsymbol{\epsilon}) \\ &= \boldsymbol{\beta}_{J_-}' \left( \frac{1}{n} \mathbf{X}_{J_-}' (\mathbf{I} - \mathbf{M}_J) \mathbf{X}_{J_-} \right) \boldsymbol{\beta}_{J_-} + \frac{1}{n} \boldsymbol{\epsilon}' \boldsymbol{\epsilon} - \frac{1}{n} \boldsymbol{\epsilon}' \mathbf{M}_J \boldsymbol{\epsilon} + \frac{2}{n} \boldsymbol{\beta}_{J_-}' \mathbf{X}_{J_-}' (\mathbf{I} - \mathbf{M}_J) \boldsymbol{\epsilon} \\ &\equiv \text{(I)} + \text{(II)} + \text{(III)} + \text{(IV)}. \end{aligned} \tag{33}$$

Define  $\hat{B} = (\mathbf{X}'_J \mathbf{X}_J)^{-1} \mathbf{X}'_J \mathbf{X}_{J-}$ . Then

$$\begin{aligned} & \mathbf{X}'_{J-} (\mathbf{I} - \mathbf{M}_J) \mathbf{X}_{J-} \quad (\mathbf{I}: \text{this is the identity matrix of dimension } \#(J_-)) \\ = & (\mathbf{I}, -\hat{B}') \begin{pmatrix} \mathbf{X}'_{J-} \\ \mathbf{X}_J \end{pmatrix} (\mathbf{X}_{J-}, \mathbf{X}_J) \begin{pmatrix} \mathbf{I} \\ -\hat{B} \end{pmatrix} \\ = & (\mathbf{I}, -\hat{B}') \mathbf{X}'_{J*} \mathbf{X}_{J*} \begin{pmatrix} \mathbf{I} \\ -\hat{B} \end{pmatrix}, \end{aligned}$$

and hence

$$\begin{aligned} \lambda_{\min}(\mathbf{X}'_{J-} (\mathbf{I} - \mathbf{M}_J) \mathbf{X}_{J-}) & \geq \lambda_{\min}(\mathbf{X}'_{J*} \mathbf{X}_{J*}) \lambda_{\min} \left( (\mathbf{I}, -\hat{B}') \begin{pmatrix} \mathbf{I} \\ -\hat{B} \end{pmatrix} \right) \\ & = \lambda_{\min}(\mathbf{X}'_{J*} \mathbf{X}_{J*}) \lambda_{\min}(\mathbf{I} + \hat{B}' \hat{B}) \\ & \geq \lambda_{\min}(\mathbf{X}'_{J*} \mathbf{X}_{J*}). \end{aligned}$$

Assume

$$\frac{1}{n} \mathbf{X}' \mathbf{X} \xrightarrow{n \rightarrow \infty} \mathbf{R} \text{ with } \lambda_{\min}(\mathbf{R}) > 0. \quad (34)$$

Then,

$$\begin{aligned} \lambda_{\min} \left( \frac{\mathbf{X}'_{J-} (\mathbf{I} - \mathbf{M}_J) \mathbf{X}_{J-}}{n} \right) &\geq \lambda_{\min} \left( \frac{\mathbf{X}'_{J*} \mathbf{X}_{J*}}{n} \right) \\ &\stackrel{\text{this is obvious}}{\geq} \lambda_{\min} \left( \frac{\mathbf{X}' \mathbf{X}}{n} \right) \\ &\stackrel{\text{by assumption}}{\longrightarrow} \lambda_{\min}(\mathbf{R}) > 0. \end{aligned} \quad (35)$$

- Therefore,

$$(I) \xrightarrow{\text{by (34)}} V(J_-) \geq \|\beta_{J_-}\|^2 \lim_{n \rightarrow \infty} \lambda_{\min} \left( \frac{1}{n} \mathbf{X}'_{J_-} (\mathbf{I} - \mathbf{M}_J) \mathbf{X}_{J_-} \right) > 0, \quad (36)$$

where

$$V(J_-) = \beta'_{J_-} \left( \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}'_{J_-} (\mathbf{I} - \mathbf{M}_J) \mathbf{X}_{J_-} \right) \beta_{J_-}.$$

- $\|\beta_{J_-}\|^2 > 0$  because  $\beta_{J_-}$  contains non-zero coefficients
- by (35),  $\lim_{n \rightarrow \infty} \lambda_{\min} \left( \frac{1}{n} \mathbf{X}'_{J_-} (\mathbf{I} - \mathbf{M}_J) \mathbf{X}_{J_-} \right) > 0$
- Moreover, it is easy to see that

$$\frac{1}{n} \epsilon' \epsilon \rightarrow \sigma^2 \text{ in probability,} \quad (37)$$

$$\frac{1}{n} \epsilon' \mathbf{M}_J \epsilon \rightarrow 0 \text{ in probability,} \quad (38)$$

and

$$E(\text{IV})^2 \leq \frac{4\sigma^2}{n} \|\beta_{J_-}\|^2 \lambda_{\max} \left( \frac{\mathbf{X}' \mathbf{X}}{n} \right) \xrightarrow{\text{by (34)}} 0. \quad (39)$$



- Now, it follows from (33) and (36)–(39) that

$$\log \hat{\sigma}_J^2 \xrightarrow{pr.} \log(\sigma^2 + V(J_-)),$$

which, together with the second equation of (32), yields

$$\log \hat{\sigma}_J^2 - \log \hat{\sigma}_{J^*}^2 \xrightarrow{pr.} \log(\sigma^2 + V(J_-)) - \log \sigma^2 > 0. \quad (40)$$

- Equation (40) and the first equation of (32) imply

$$\begin{aligned} & P(\text{BIC}(J) \leq \text{BIC}(J^*)) \\ &= P\left(\log \hat{\sigma}_J^2 - \log \hat{\sigma}_{J^*}^2 \leq \frac{\log n}{n}(\sharp(J^*) - \sharp(J))\right) \xrightarrow{n \rightarrow \infty} 0, \quad \left(\frac{\log n}{n} \rightarrow 0\right) \end{aligned}$$

which is (28). Thus, the proof is complete.

# A Newton-Raphson method and its asymptotics

- Define

$$\hat{\boldsymbol{\eta}}^{\text{New}} = \tilde{\boldsymbol{\eta}} - \left( \frac{\partial^2 \ell(\tilde{\boldsymbol{\eta}})}{\partial \eta_i \partial \eta_j} \right)^{-1} \frac{\partial \ell(\tilde{\boldsymbol{\eta}})}{\partial \boldsymbol{\eta}},$$

where  $\ell(\boldsymbol{\eta})$  is the log-likelihood function and  $\tilde{\boldsymbol{\eta}}$  is an initial estimate of the true parameter  $\boldsymbol{\eta}_0$ .

- In the following, I'll show that if

$$\|\tilde{\boldsymbol{\eta}} - \boldsymbol{\eta}_0\| = o_p(n^{-\frac{1}{4}}), \quad (41)$$

then

$$\sqrt{n}(\hat{\boldsymbol{\eta}}^{\text{New}} - \boldsymbol{\eta}_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{I}^{-1}(\boldsymbol{\eta}_0)),$$

where  $\mathbf{I}(\boldsymbol{\eta}_0)$  is the Fisher information matrix.

# Sketch

$$\hat{\boldsymbol{\eta}}^{\text{New}} = \tilde{\boldsymbol{\eta}} - \left( \frac{\partial^2 \ell(\tilde{\boldsymbol{\eta}})}{\partial \eta_i \partial \eta_j} \right)^{-1} \left( \frac{\partial \ell(\tilde{\boldsymbol{\eta}})}{\partial \boldsymbol{\eta}} - \frac{\partial \ell(\boldsymbol{\eta}_0)}{\partial \boldsymbol{\eta}} \right) - \left( \frac{\partial^2 \ell(\tilde{\boldsymbol{\eta}})}{\partial \eta_i \partial \eta_j} \right)^{-1} \frac{\partial \ell(\boldsymbol{\eta}_0)}{\partial \boldsymbol{\eta}}$$

By Taylor's theorem,

$$\frac{\partial \ell(\boldsymbol{\eta}_0)}{\partial \boldsymbol{\eta}} = \frac{\partial \ell(\tilde{\boldsymbol{\eta}})}{\partial \boldsymbol{\eta}} + \frac{\partial^2 \ell(\tilde{\boldsymbol{\eta}})}{\partial \eta_i \partial \eta_j} (\boldsymbol{\eta}_0 - \tilde{\boldsymbol{\eta}}) + \begin{pmatrix} \gamma_1 \\ \vdots \\ \gamma_p \end{pmatrix}, \quad (42)$$

where

$$\gamma_k = \frac{1}{2} \sum_{i=1}^p \sum_{j=1}^p \frac{\partial^3 \ell(\boldsymbol{\eta}_k^*)}{\partial \eta_k \partial \eta_i \partial \eta_j} (\eta_{0,i} - \tilde{\eta}_i)(\eta_{0,j} - \tilde{\eta}_j)$$

with

$$\begin{pmatrix} \tilde{\eta}_1 \\ \vdots \\ \tilde{\eta}_p \end{pmatrix} = \tilde{\boldsymbol{\eta}}, \quad \begin{pmatrix} \eta_{0,1} \\ \vdots \\ \eta_{0,p} \end{pmatrix} = \boldsymbol{\eta}_0 \quad \text{and} \quad \|\tilde{\boldsymbol{\eta}} - \boldsymbol{\eta}_0\| \geq \|\boldsymbol{\eta}_k^* - \boldsymbol{\eta}_0\|, k = 1, \dots, p.$$

(assuming there are  $p$  parameters)

- Assume Condition (ii) in the proof of the consistency of the MLE holds, and there exists  $\delta > 0$  such that

$$\sup_{\boldsymbol{\eta} \in B_\delta(\boldsymbol{\eta}_0)} \frac{1}{n} \frac{\partial^3 \ell(\boldsymbol{\eta})}{\partial \eta_k \partial \eta_i \partial \eta_j} = O_p(1).$$

- Then, by (41) and (42), one gets

$$\hat{\boldsymbol{\eta}}^{\text{New}} = \tilde{\boldsymbol{\eta}} - (\tilde{\boldsymbol{\eta}} - \boldsymbol{\eta}_0) + o_p(n^{-\frac{1}{2}}) - \left( \frac{\partial^2 \ell(\tilde{\boldsymbol{\eta}})}{\partial \eta_i \partial \eta_j} \right)^{-1} \frac{\partial \ell(\boldsymbol{\eta}_0)}{\partial \boldsymbol{\eta}},$$

and hence

$$\sqrt{n}(\hat{\boldsymbol{\eta}}^{\text{New}} - \boldsymbol{\eta}_0) = o_p(1) - \left( \frac{1}{n} \frac{\partial^2 \ell(\tilde{\boldsymbol{\eta}})}{\partial \eta_i \partial \eta_j} \right)^{-1} \frac{1}{\sqrt{n}} \frac{\partial \ell(\boldsymbol{\eta}_0)}{\partial \boldsymbol{\eta}}. \quad (43)$$

By (41) and an argument similar to that used to prove (42), the RHS of (43) becomes

$$\left( -\frac{1}{n} \frac{\partial^2 \ell(\boldsymbol{\eta}_0)}{\partial \eta_i \partial \eta_j} \right)^{-1} \frac{1}{\sqrt{n}} \frac{\partial \ell(\boldsymbol{\eta}_0)}{\partial \boldsymbol{\eta}} + o_p(1),$$

and hence the desired conclusion follows from

$$\frac{1}{\sqrt{n}} \frac{\partial \ell(\boldsymbol{\eta}_0)}{\partial \boldsymbol{\eta}} \xrightarrow{d} N(\mathbf{0}, \mathbf{I}(\boldsymbol{\eta}_0)),$$

and

$$-\frac{1}{n} \frac{\partial^2 \ell(\boldsymbol{\eta}_0)}{\partial \eta_i \partial \eta_j} \xrightarrow{pr.} \mathbf{I}(\boldsymbol{\eta}_0).$$

### Question

Please show that if (41) is replaced by

$$\|\tilde{\boldsymbol{\eta}} - \boldsymbol{\eta}_0\| = O_p(n^{-\frac{1}{q}})$$

for some  $q \geq 4$ , then

$$\|\hat{\boldsymbol{\eta}}^{\text{New}} - \boldsymbol{\eta}_0\| = O_p(n^{-\frac{2}{q}}).$$

# Example

- Assume

$$y_i = \mathbf{x}_i' \boldsymbol{\beta}_0 + a_i, \quad i = 1, \dots, n,$$

where  $a_i = \sigma_i \epsilon_i$ ,  $\epsilon_i \stackrel{\text{iid}}{\sim} N(0, 1)$ ,  $\sigma_i^2 = \exp\{\mathbf{x}_i' \boldsymbol{\alpha}_0\}$ , and  $\mathbf{x}_i$  is a  $(p+1)$ -dimensional explanatory vector with 1 as its first component (namely the intercept term is included).

- Assume also  $1 \leq \|\mathbf{x}_i\| \leq M$  for some  $1 < M < \infty$  and

$$\frac{1}{n} \mathbf{X}' \mathbf{X} \rightarrow \mathbf{R} \text{ (p.d.)}, \quad \text{where } \mathbf{X} = \begin{pmatrix} \mathbf{x}_1' \\ \vdots \\ \mathbf{x}_n' \end{pmatrix}.$$

Then, it is not difficult to see that the LSE of  $\boldsymbol{\beta}_0$ ,

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y}, \quad \text{where } \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix},$$

satisfies

$$\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0 = O_p(n^{-\frac{1}{2}}).$$

(44)

- In the following, I'll construct an initial estimate of

$$\theta_0 = \begin{pmatrix} u + \alpha_{0,0} \\ \alpha_{0,1} \\ \vdots \\ \alpha_{0,p} \end{pmatrix}, \text{ where } \begin{pmatrix} \alpha_{0,0} \\ \vdots \\ \alpha_{0,p} \end{pmatrix} = \alpha_0 \text{ and } u = E(\log \epsilon_i^2).$$

- Define

$$\tilde{a}_i = y_i - \mathbf{x}'_i \hat{\beta} = a_i - \mathbf{x}'_i (\hat{\beta} - \beta_0),$$

and

$$\hat{a}_i = \begin{cases} \tilde{a}_i, & \text{if } |\tilde{a}_i| > n^{-\xi}, \\ n^{-\xi}, & \text{if } |\tilde{a}_i| \leq n^{-\xi}, \end{cases}$$

where  $0 < \xi < \frac{1}{2} - \theta_1$ . Let  $A_n = \{\|\hat{\beta} - \beta_0\| \leq n^{-\frac{1}{2} + \theta_1}\}$ ,  $\theta_1 > 0$  is a small positive number.

- By (44), it holds that

$$\lim_{n \rightarrow \infty} P(A_n) = 1. \quad (45)$$

- Let

$$\hat{\boldsymbol{\theta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{z}, \text{ where } \mathbf{z} = \begin{pmatrix} \log \hat{a}_1^2 \\ \vdots \\ \log \hat{a}_n^2 \end{pmatrix}.$$

- Since

$$\log \hat{a}_i^2 = \mathbf{x}_i'\boldsymbol{\theta}_0 + (\log \hat{a}_i^2 - \log a_i^2) + (\log \epsilon_i^2 - u),$$

we obtain

$$\begin{aligned} \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 &= \left(\frac{1}{n}\mathbf{X}'\mathbf{X}\right)^{-1} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i (\log \epsilon_i^2 - u) \\ &\quad + \left(\frac{1}{n}\mathbf{X}'\mathbf{X}\right)^{-1} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i (\log \hat{a}_i^2 - \log a_i^2) I_{|a_i| < n^{-\theta}} \\ &\quad + \left(\frac{1}{n}\mathbf{X}'\mathbf{X}\right)^{-1} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i (\log \hat{a}_i^2 - \log a_i^2) I_{|a_i| > n^{-\theta}} \\ &\equiv \text{(I)} + \text{(II)} + \text{(III)}, \end{aligned}$$

where  $\xi < \theta < \frac{1}{2} - \theta_1$ .



- It is easy to see that

$$(I) = O_p(n^{-\frac{1}{2}}). \quad (46)$$

- Denote  $\mathbf{x}_i$  by  $\begin{pmatrix} x_{i1} \\ \vdots \\ x_{i,p+1} \end{pmatrix}$ . Then, we have for  $1 \leq j \leq p+1$ ,

$$\begin{aligned} & E \left\{ \frac{1}{n} \sum_{i=1}^n |x_{ij}(\log \hat{a}_i^2 - \log a_i^2)| I_{|a_i| < n^{-\theta}} I_{A_n} \right\} \\ & \stackrel{\text{why?}}{\leq} C \log n P(a_i^2 < n^{-2\theta}) + CP^{\frac{1}{2}}(a_i^2 < n^{-2\theta}) \quad (C : \text{some positive constant}) \\ & \stackrel{\text{why?}}{\leq} C^* n^{-\frac{\theta}{2}}. \quad (C^* : \text{some positive constant}) \end{aligned}$$

This, the positive definiteness of  $\mathbf{R}$ , and (45) yield

$$(II) = O_p(n^{-\frac{\theta}{2}}). \quad (47)$$

Moreover,

$$\begin{aligned}
 & E \left| \frac{1}{n} \sum_{i=1}^n x_{ij} (\log \hat{a}_i^2 - \log a_i^2) I_{|a_i| > n^{-\theta}} I_{A_n} \right| \\
 \stackrel{\text{why?}}{\leq} & E \left| \frac{C}{n} \sum_{i=1}^n |x_{ij}| \frac{|\mathbf{x}'_i(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)|}{|a_i|} I_{|a_i| > n^{-\theta}} I_{A_n} \right| \quad (C : \text{some positive constant}) \\
 \stackrel{\text{why?}}{\leq} & C^* n^{\theta - \frac{1}{2} + \theta_1}. \quad (C^* : \text{some positive constant})
 \end{aligned}$$

This, the positive definiteness of  $\mathbf{R}$ , and (45) imply

$$(III) = O_p(n^{-\frac{1}{2} + \theta_1 + \theta}),$$

which, together with (46) and (47), gives

$$\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 = O_p(n^{-\zeta}) \quad (48)$$

with  $\zeta = \min\{\frac{1}{2} - \theta_1 - \theta, \frac{\theta}{2}\}$ .

## Remark

- Recall that the location-dispersion model is

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + \exp \left\{ \frac{1}{2} \mathbf{x}_i' \boldsymbol{\alpha} \right\} \epsilon_i, \quad i = 1, \dots, n,$$

with  $\epsilon_i \stackrel{\text{i.i.d.}}{\sim} N(0, 1)$  and the log-likelihood function is

$$\ell(\boldsymbol{\beta}, \boldsymbol{\alpha}) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^n \mathbf{x}_i' \boldsymbol{\alpha} - \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i' \boldsymbol{\beta})^2 e^{-\mathbf{x}_i' \boldsymbol{\alpha}}.$$

- The function  $-\ell(\boldsymbol{\beta}, \boldsymbol{\alpha})$  is not jointly convex in coefficients  $(\boldsymbol{\beta}, \boldsymbol{\alpha})$ . This also reveals the importance of finding a good initial estimate of the true parameters.

## Model

$$\begin{aligned}y_t &= \beta x_t + \eta_t, \quad t = 1, 2, \dots, n, \\ \eta_t &= \sigma_t \epsilon_t, \\ \sigma_t^2 &= \exp\{\alpha x_t\}\end{aligned}$$

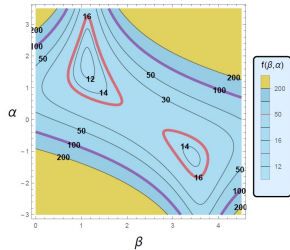
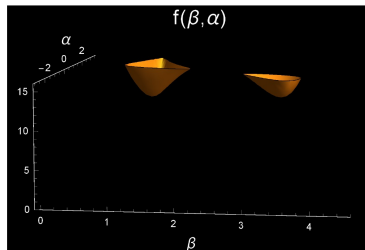
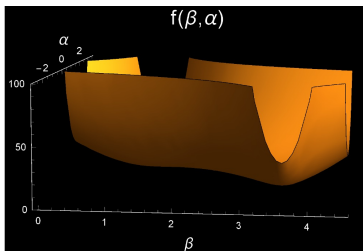
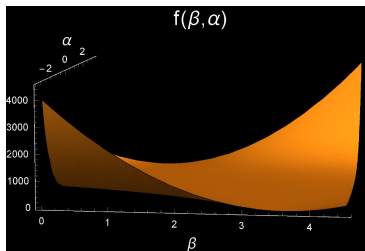
## Function

$$f(\beta, \alpha) = \sum_{i=1}^n (\alpha x_i) + \sum_{j=1}^n \frac{(y_j - \beta x_j)^2}{\exp\{\alpha x_j\}}$$

## Setting

- $x_t \stackrel{\text{i.i.d.}}{\sim} N(0, 1)$
- $\epsilon_t \stackrel{\text{i.i.d.}}{\sim} N(0, 1)$
- $\beta = 1$  and  $\alpha = 2$
- $n = 10$

# 3D plot and contour plot



**Question**

Define

$$\tilde{\alpha}_0 = \log \frac{1}{n} \sum_{i=1}^n \frac{\hat{a}_i^2}{e^{x_{i1}\hat{\theta}_1 + \dots + x_{ip}\hat{\theta}_p}},$$

where

$$\begin{pmatrix} \hat{\theta}_0 \\ \hat{\theta}_1 \\ \vdots \\ \hat{\theta}_p \end{pmatrix} = \hat{\boldsymbol{\theta}}.$$

- a. Please show that  $\tilde{\alpha}_0 - \alpha_{0,0} = O_p(n^{-\delta})$  for some  $\delta > 0$ .
- b. Can you obtain similar results when  $\epsilon_i$ 's are non-Gaussian?

# Asymptotics for the MLEs of Weibull distribution parameters

- Weibull distribution

$$f(x) = \frac{\alpha_0}{\lambda_0^{\alpha_0}} x^{\alpha_0-1} \exp \left\{ - \left( \frac{x}{\lambda_0} \right)^{\alpha_0} \right\}, \quad (49)$$

where  $x > 0$ ,  $\lambda_0 > 0$  is called the scale parameter, and  $\alpha_0 > 0$  is called the shape parameter.

- The density function (49) can be rewritten as

$$f(x) = \frac{\alpha_0}{\eta_0} x^{\alpha_0-1} \exp \left\{ - \frac{x^{\alpha_0}}{\eta_0} \right\}, \quad x > 0, \alpha_0 > 0, \eta_0 > 0,$$

in which  $\eta_0 = \lambda_0^{\alpha_0}$ .

Define the average log-likelihood function

$$\frac{1}{n}\ell(\alpha, \eta) = \log \frac{\alpha}{\eta} + (\alpha - 1)\frac{1}{n} \sum_{i=1}^n \log x_i - \frac{1}{\eta} \frac{1}{n} \sum_{i=1}^n x_i^{\alpha}.$$

Then

$$\begin{aligned} & \frac{1}{n} \{ \ell(\alpha_0, \eta_0) - \ell(\alpha, \eta) \} \\ = & \log \frac{\alpha_0}{\eta_0} - \log \frac{\alpha}{\eta} + (\alpha_0 - \alpha) \frac{1}{n} \sum_{i=1}^n (\log x_i - E(\log x_i)) \\ & + (\alpha_0 - \alpha) E(\log x_i) - \frac{1}{n} \sum_{i=1}^n \frac{x_i^{\alpha_0} - E x_i^{\alpha_0}}{\eta_0} \\ & + \frac{1}{n} \sum_{i=1}^n \frac{x_i^{\alpha} - E x_i^{\alpha}}{\eta} + \frac{E x_i^{\alpha}}{\eta} - \frac{E x_i^{\alpha_0}}{\eta_0}. \end{aligned} \tag{50}$$



- Note that

$$\begin{aligned}
 E(\log x_i) &= \frac{1}{\alpha_0} \int_0^\infty \left( \log \frac{x_i^{\alpha_0}}{\eta_0} + \log \eta_0 \right) \frac{\alpha_0}{\eta_0} x_i^{\alpha_0-1} e^{-\frac{x_i^{\alpha_0}}{\eta_0}} dx_i \\
 &= \frac{1}{\alpha_0} \int_0^\infty \log u e^{-u} du + \frac{1}{\alpha_0} \log \eta_0 \\
 &= \frac{1}{\alpha_0} (-\gamma + \log \eta_0),
 \end{aligned}$$

where  $\gamma = \lim_{n \rightarrow \infty} (\log n - \sum_{i=1}^n \frac{1}{i}) \sim 0.5772$ . ( $\gamma$ : Euler-Mascheroni constant)

- In addition, it holds that for  $m > -\alpha_0$ ,

$$E(x_i^m) = \eta_0^{\frac{m}{\alpha_0}} \Gamma\left(\frac{m}{\alpha_0} + 1\right),$$

and by the Weierstrass product for the  $\Gamma$  function,

$$\frac{\Gamma'(z+1)}{\Gamma(z)} = -\gamma + \sum_{i \geq 1} \frac{1}{i} - \frac{1}{i+z}, \quad z > 0. \quad (51)$$

- (51) implies

$$\frac{d}{dz} \frac{\Gamma'(z+1)}{\Gamma(z)} = \frac{\Gamma(z+1)\Gamma''(z+1) - (\Gamma'(z+1))^2}{\Gamma^2(z)} = \sum_{i \geq 1} \frac{1}{(i+z)^2} > 0, \quad z > 0,$$

and hence

$$\frac{\Gamma''(z+1)}{\Gamma(z)} \geq \left( \frac{\Gamma'(z+1)}{\Gamma(z)} \right)^2, \quad z > 0. \quad (52)$$

- Now, for  $(\alpha, \eta) \in [\delta_1, M_1] \times [\delta_2, M_2]$ , where  $0 < \delta_1, \delta_2 < \infty$  are small constants and  $0 < M_1, M_2 < \infty$  are large constants,

$$\begin{aligned} & \frac{1}{n} \{ \ell(\alpha_0, \eta_0) - \ell(\alpha, \eta) \} \\ &= (\alpha_0 - \alpha) \frac{1}{n} \sum_{i=1}^n \{ \log x_i - E(\log x_i) \} - \frac{1}{n} \sum_{i=1}^n \left( \frac{x_i^{\alpha_0}}{\eta_0} - 1 \right) \\ & \quad + \frac{1}{n} \sum_{i=1}^n \left\{ \frac{x_i^\alpha}{\eta} - \frac{E x_i^\alpha}{\eta} \right\} + g(\alpha_0, \eta_0) - g(\alpha, \eta) \\ &\equiv \text{(I)} + \text{(II)} + \text{(III)} + \text{(IV)}, \end{aligned} \quad (53)$$

where  $\text{(IV)} = g(\alpha_0, \eta_0) - g(\alpha, \eta)$  and

$$g(\alpha, \eta) = \log \frac{\alpha}{\eta} + \alpha E(\log x_i) - \frac{\eta_0^{\frac{\alpha}{\alpha_0}} \Gamma(\frac{\alpha}{\alpha_0} + 1)}{\eta}.$$

In the following, I shall show that

(i)  $(\alpha_0, \eta_0)$  is the only critical point satisfying

$$\frac{\partial g(\alpha, \eta)}{\partial \alpha} = 0 \quad \text{and} \quad \frac{\partial g(\alpha, \eta)}{\partial \eta} = 0,$$

(ii) there exist small positive constants  $s_1$  and  $s_2$  such that

$$g(\alpha_0, \eta_0) - g(\alpha, \eta) \geq s_2 \|\mathbf{v} - \mathbf{v}_0\|^2$$

for all  $\mathbf{v} = (\alpha, \eta)' \in B_{s_1}(\mathbf{v}_0)$ , where  $\mathbf{v}_0 = (\alpha_0, \eta_0)'$ ,

(iii)  $(\alpha_0, \eta_0)$  is the unique maximizer of  $g(\alpha, \eta)$ .

To show (i), note that

$$\frac{\partial g(\alpha, \eta)}{\partial \eta} = 0 \iff \eta = \eta_0^{\frac{\alpha}{\alpha_0}} \Gamma \left( \frac{\alpha}{\alpha_0} + 1 \right),$$

and

$$\frac{d}{d\alpha} g \left( \alpha, \eta_0^{\frac{\alpha}{\alpha_0}} \Gamma \left( \frac{\alpha}{\alpha_0} + 1 \right) \right) = 0 \overset{\text{why?}}{\iff} \sum_{i \geq 1} \left( \frac{1}{i} - \frac{1}{i + \frac{\alpha}{\alpha_0}} \right) = \frac{1}{\frac{\alpha}{\alpha_0}}. \quad (54)$$

- Since the RHS (LHS) of (54) is decreasing (increasing) in  $\frac{\alpha}{\alpha_0}$ , and  $\frac{\alpha}{\alpha_0} = 1$  is a solution of the equation, the desired conclusion (i) follows (why?).
- To show (ii), note first that by (i) and Taylor's theorem, one has for  $\|\mathbf{v} - \mathbf{v}_0\| \leq \delta$  with  $\delta$  being arbitrarily small,

$$g(\alpha, \eta) = g(\alpha_0, \eta_0) - \frac{1}{2}(\alpha - \alpha_0, \eta - \eta_0)' \begin{pmatrix} \frac{\partial^2 g(\alpha^*, \eta^*)}{\partial \alpha^2} & \frac{\partial^2 g(\alpha^*, \eta^*)}{\partial \alpha \partial \eta} \\ \frac{\partial^2 g(\alpha^*, \eta^*)}{\partial \eta \partial \alpha} & \frac{\partial^2 g(\alpha^*, \eta^*)}{\partial \eta^2} \end{pmatrix} \begin{pmatrix} \alpha - \alpha_0 \\ \eta - \eta_0 \end{pmatrix}, \quad (55)$$

where

$$g(\alpha, \eta) = \frac{\eta_0^{\frac{\alpha}{\alpha_0}} \Gamma\left(\frac{\alpha}{\alpha_0} + 1\right)}{\eta} - \log \frac{\alpha}{\eta} \quad \text{and} \quad \left\| \begin{pmatrix} \alpha^* - \alpha_0 \\ \eta^* - \eta_0 \end{pmatrix} \right\| \leq \left\| \begin{pmatrix} \alpha - \alpha_0 \\ \eta - \eta_0 \end{pmatrix} \right\| \leq \delta.$$

It is straightforward to see that

$$\frac{\partial^2 g(\alpha, \eta)}{\partial \alpha^2} = \frac{\eta_0^{\frac{\alpha}{\alpha_0}}}{\eta} \left[ \left( \frac{\log \eta_0}{\alpha_0} \right)^2 \Gamma \left( \frac{\alpha}{\alpha_0} + 1 \right) + \frac{2 \log \eta_0}{\alpha_0^2} \Gamma' \left( \frac{\alpha}{\alpha_0} + 1 \right) + \frac{\Gamma'' \left( \frac{\alpha}{\alpha_0} + 1 \right)}{\alpha_0^2} \right] + \frac{1}{\alpha^2},$$

$$\frac{\partial^2 g(\alpha, \eta)}{\partial \eta^2} = \frac{\eta_0^{\frac{\alpha}{\alpha_0}}}{\eta} \left( \frac{2 \Gamma \left( \frac{\alpha}{\alpha_0} + 1 \right)}{\eta^2} \right) - \frac{1}{\eta^2},$$

$$\frac{\partial^2 g(\alpha, \eta)}{\partial \alpha \partial \eta} = \frac{\eta_0^{\frac{\alpha}{\alpha_0}}}{\eta} \left\{ -\frac{1}{\alpha_0 \eta} \left[ \log \eta_0 \Gamma \left( \frac{\alpha}{\alpha_0} + 1 \right) + \Gamma' \left( \frac{\alpha}{\alpha_0} + 1 \right) \right] \right\}.$$

Since  $\|\mathbf{v} - \mathbf{v}_0\| \leq \delta$  with  $\delta$  being arbitrarily small, there exists an arbitrarily small  $\theta_\delta > 0$  such that for all  $\|\mathbf{v} - \mathbf{v}_0\| \leq \delta$ ,

$$\frac{\eta_0^{\frac{\alpha}{\alpha_0}}}{\eta} \left( \frac{2\Gamma\left(\frac{\alpha}{\alpha_0} + 1\right)}{\eta^2} \right) - \frac{1}{\eta^2} \geq (1 - \theta_\delta) \frac{\eta_0^{\frac{\alpha}{\alpha_0}}}{\eta} \frac{\Gamma\left(\frac{\alpha}{\alpha_0} + 1\right)}{\eta^2},$$

which, together with (52), yields

$$\begin{aligned} & \det \begin{pmatrix} \frac{\partial^2 g(\alpha, \eta)}{\partial \alpha^2} & \frac{\partial^2 g(\alpha, \eta)}{\partial \alpha \partial \eta} \\ \frac{\partial^2 g(\alpha, \eta)}{\partial \eta \partial \alpha} & \frac{\partial^2 g(\alpha, \eta)}{\partial \eta^2} \end{pmatrix} \\ & \geq \left( \frac{\eta_0^{\frac{\alpha}{\alpha_0}}}{\eta} \right)^2 \left( \frac{1}{\alpha_0 \eta} \right)^2 \left\{ (1 - \theta_\delta) \left[ (\log \eta_0)^2 \Gamma^2 \left( \frac{\alpha}{\alpha_0} + 1 \right) \right. \right. \\ & \quad \left. \left. + 2 \log \eta_0 \Gamma' \left( \frac{\alpha}{\alpha_0} + 1 \right) \Gamma \left( \frac{\alpha}{\alpha_0} + 1 \right) + \Gamma'' \left( \frac{\alpha}{\alpha_0} + 1 \right) \Gamma \left( \frac{\alpha}{\alpha_0} + 1 \right) \right] \right. \\ & \quad \left. - \left[ \left( \log \eta_0 \Gamma \left( \frac{\alpha}{\alpha_0} + 1 \right) \right)^2 + 2 \log \eta_0 \Gamma \left( \frac{\alpha}{\alpha_0} + 1 \right) \Gamma' \left( \frac{\alpha}{\alpha_0} + 1 \right) + \left( \Gamma' \left( \frac{\alpha}{\alpha_0} + 1 \right) \right)^2 \right] \right\} \\ & > \underline{c}, \end{aligned}$$

for some small  $\underline{c} > 0$  depending only on  $\delta$  and  $\theta_\delta$ .

- Therefore,

$$\inf_{\mathbf{v} \in \mathcal{B}_\delta(\mathbf{v}_0)} \lambda_{\min} \left( \begin{pmatrix} \frac{\partial^2 g(\alpha, \eta)}{\partial \alpha^2} & \frac{\partial^2 g(\alpha, \eta)}{\partial \alpha \partial \eta} \\ \frac{\partial^2 g(\alpha, \eta)}{\partial \eta \partial \alpha} & \frac{\partial^2 g(\alpha, \eta)}{\partial \eta^2} \end{pmatrix} \right) \stackrel{\text{why?}}{>} \underline{c}^*, \quad (56)$$

for some  $\underline{c}^* > 0$ .

- By (55) and (56), it holds that for  $\|\mathbf{v} - \mathbf{v}_0\| \leq \delta$ ,

$$g(\alpha_0, \eta_0) - g(\alpha, \eta) \geq \frac{\underline{c}^*}{2} \|\mathbf{v} - \mathbf{v}_0\|^2,$$

and hence (ii) follows.

- Now, (iii) follows directly from (i), (ii) and the differentiability of  $g(\alpha, \beta)$ .

- In view of (i)–(iii) and (53), the consistency of

$$(\hat{\alpha}, \hat{\eta}) = \operatorname{argmax}_{(\alpha, \eta) \in [\delta_1, M_1] \times [\delta_2, M_2]} \frac{1}{n} \ell(\alpha, \eta)$$

is ensured by

$$\sup_{\alpha \in [\delta_1, M_1]} |(\text{I})| = o_p(1), \quad |(\text{II})| = o_p(1), \quad (57)$$

$$\sup_{(\alpha, \eta) \in [\delta_1, M_1] \times [\delta_2, M_2]} |(\text{III})| = o_p(1). \quad (58)$$

- (57) is an immediate consequence of the (classical) law of large numbers, whereas (58) relies on the so-called uniform law of large numbers.
- In the following, I'll provide a proof of (58).



- Note first that (58) is guaranteed by (why?)

$$\sup_{\alpha \in [\delta_1, M_1]} \left| \frac{1}{n} \sum_{i=1}^n x_i^\alpha - E(x_i^\alpha) \right| = o_p(1). \quad (59)$$

- Define  $g_t(\alpha) = x_t^\alpha - E(x_t^\alpha)$ . Then,

$$\begin{aligned} & E \sup_{\alpha \in [\delta_1, M_1]} \left( \frac{1}{n} \sum_{t=1}^n g_t(\alpha) - g_t(\alpha_0) \right)^2 \\ & \stackrel{\text{why?}}{\leq} (M_1 - \delta_1) E \left\{ \int_{\delta_1}^{M_1} \left( \frac{1}{n} \sum_{t=1}^n g'_t(x) \right)^2 dx \right\} \\ & \stackrel{\text{why?}}{\leq} (M_1 - \delta_1)^2 \sup_{\alpha \in [\delta_1, M_1]} E \left( \frac{1}{n} \sum_{t=1}^n g'_t(\alpha) \right)^2 \\ & \stackrel{\text{why?}}{\leq} \frac{(M_1 - \delta_1)^2}{\eta} \sup_{\alpha \in [\delta_1, M_1]} \text{Var}(x_1^\alpha \log x_1) = O(n^{-1}), \end{aligned}$$

which implies

$$\sup_{\alpha \in [\delta_1, M_1]} \frac{1}{n} \sum_{t=1}^n g_t(\alpha) - g_t(\alpha_0) \stackrel{\text{why?}}{=} O_p(n^{-\frac{1}{2}}). \quad (60)$$

- Moreover, it is easy to see that

$$\frac{1}{n} \sum_{t=1}^n g_t(\alpha_0) = O_p(n^{-\frac{1}{2}}). \quad (61)$$

- Combining (60) and (61) gives (59). Thus, the proof is complete.

### Question

Please find the limiting distribution of  $(\hat{\alpha}, \hat{\eta})$ .