

$$\begin{aligned}
E[\hat{\sigma}_{MLE}^2] &= E\left[\frac{1}{n} \sum_{i=1}^n \frac{1}{4} (y_{i1} - y_{i2})^2\right] \\
(\text{i.i.d.}) &= \frac{1}{4} E[(y_{i1} - y_{i2})^2] \\
&= \frac{1}{4} E[(y_{i1} - \mu_i) - (y_{i2} - \mu_i)]^2 \\
&= \frac{1}{4} E[(y_{i1} - \mu_i)^2 - 2(y_{i1} - \mu_i)(y_{i2} - \mu_i) + (y_{i2} - \mu_i)^2] \\
&= \frac{1}{4} [\sigma^2 - 0 + \sigma^2] \\
&= \frac{\sigma^2}{2} \neq \sigma^2!
\end{aligned}$$

So as $n \rightarrow \infty$, $\hat{\sigma}_{MLE}^2 \xrightarrow{p} \frac{\sigma^2}{2}$ by WLLN
This seems bad.

Happens because # of nuisance parameters grows w/ n (# groups).
A reworking of the data seems more promising. Let,

$$V_i = \frac{Y_{i1} - Y_{i2}}{\sqrt{2}} \quad \text{and} \quad W_i = \frac{Y_{i1} + Y_{i2}}{\sqrt{2}}$$

Because Y_{ij} are Gaussian,

$$V_i \sim N(0, \sigma^2) \quad \text{and} \quad W_i \sim N(\sqrt{2}\mu_i, \sigma^2) \quad \text{Also, } V_i \perp W_i!$$

$$\begin{bmatrix} V_i \\ W_i \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} Y_{i1} \\ Y_{i2} \end{bmatrix} \Rightarrow \text{Var}\left(\begin{bmatrix} V_i \\ W_i \end{bmatrix}\right) = \frac{1}{2} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} = \sigma^2 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

consider the density of $V \in W$:

$$f_{VW}(\underline{v}, \underline{w}; \sigma^2, \mu_1, \dots, \mu_n) = \underbrace{f_V(\underline{v}; \sigma^2)}_{\text{no nuisance parameters!}} f_W(\underline{w}; \mu_1, \dots, \mu_n, \sigma^2)$$

$$\Rightarrow \ell(\sigma | \underline{V}) = -n \log \sqrt{2\pi} - n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n V_i^2$$

$$\frac{\partial \ell(\sigma | \underline{V})}{\partial \sigma} = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n V_i^2 \Rightarrow \hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n V_i^2$$

A marginal likelihood approach is simple provided you can find a statistic V whose ds is free of the nuisance parameter!

For conditional likelihoods, we can often exploit the existence of sufficient statistics for the nuisance parameters under the assumption that the parameter of interest is known.

Let T_i be sufficient for a nuisance parameter (μ_i)

Then the conditional dsr of the data given $\underline{T} = (T_1, \dots, T_n)$ doesn't depend on the nuisance parameters!

\Rightarrow we can look for the conditional dsr of the data \underline{Y} .

Example (Exponential Families): The structure of exponential families is such that it is often possible to exploit their properties to eliminate nuisance parameters. Let \underline{Y} have a density of the form

$$f(\underline{y}; \underline{\eta}) = h(\underline{y}) \exp \left\{ \sum_{i=1}^s \eta_i T_i(\underline{y}) - A(\underline{\eta}) \right\},$$

then

if $\underline{\eta} = (\theta_1^T, \theta_2^T)$

$$f(\underline{y}; \underline{\theta}_1, \underline{\theta}_2) = h(\underline{y}) \exp \left\{ \sum \theta_{1i} w_i + \sum \theta_{2j} v_j - A(\underline{\theta}_1, \underline{\theta}_2) \right\}$$

These are sufficient statistics

where
the conditional distribution of \underline{W} given \underline{V} is exponential family of the form

$$f(\underline{w} | \underline{v}; \underline{\theta}_1) = q_{\underline{v}}(\underline{w}) \exp \left(\sum_{i=1}^r \theta_{1i} w_i - A_{\underline{v}}(\underline{\theta}_1) \right)$$

depends only on $\underline{\theta}_1$!

Thus, exponential families often provide an automatic procedure for finding \underline{W} and \underline{V} .

Example (Logistic Regression): For binary Y_i , the standard logistics regression model is

$$P(Y_i = 1) = p_i(\mathbf{x}_i, \boldsymbol{\beta}) = \frac{\exp(\mathbf{x}_i^\top \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})}$$

and the likelihood is

$$\begin{aligned} L(\boldsymbol{\beta} | \mathbf{Y}, \mathbf{X}) &= \prod_{i=1}^n p_i(\mathbf{x}_i, \boldsymbol{\beta})^{y_i} \{1 - p_i(\mathbf{x}_i, \boldsymbol{\beta})\}^{1-y_i} \\ &= \prod_{i=1}^n \left\{ \frac{\exp(\mathbf{x}_i^\top \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})} \right\}^{y_i} \left\{ \frac{1}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})} \right\}^{1-y_i} \\ &= \frac{\exp\left(\sum_{i=1}^n y_i (\mathbf{x}_i^\top \boldsymbol{\beta})\right)}{\prod_{i=1}^n (1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta}))} \\ &= c(\mathbf{x}, \boldsymbol{\beta}) \exp\left(\sum_{j=1}^p \beta_j \sum_{i=1}^n x_{ij} y_i\right) \end{aligned}$$

$\Rightarrow T_j = \sum_{i=1}^n x_{ij} y_i \quad j=1, \dots, p$ are sufficient for β_j 's.

Suppose $\theta_1 = \beta_k$ is the parameter of interest (e.g. treatment variable) and others are nuisance parameters.

$$\Rightarrow W_i = T_k = \sum_{i=1}^n x_{ik} y_i \quad \text{and} \quad \underline{V} = (T_1, \dots, T_{k-1}, T_{k+1}, \dots, T_p)^\top$$

joint
marginal

and the conditional density $P(W = t_k | \underline{V} = \underline{v}) = P(T_k = t_k | T_1 = t_1, \dots, T_{k-1} = t_{k-1}, T_{k+1} = t_{k+1}, \dots, T_p = t_p)$

$$\begin{aligned} &\vdots \\ &= \frac{c(t_1, \dots, t_p) \exp(\beta_k t_k)}{\sum_{\mathbf{u}} c(t_1, \dots, t_{k-1}, u, t_{k+1}, \dots, t_p) \exp(\beta_k u)} \end{aligned}$$

depends only on β_k !

1.5 The Maximum Likelihood Estimator and the Information Matrix

We have now talked about how to construct likelihoods in a variety of settings, now we can use those constructions to formalize how we make inferences about model parameters.

→ parameter estimation, hypothesis tests, confidence intervals.

We will often restrict attention to likelihoods that are continuously differentiable wrt $\underline{\theta}$.

In this case, Recall the score function

$$S(\underline{\theta}) = S(\mathbf{Y}, \underline{\theta}) = \begin{pmatrix} \frac{\partial \ell(\underline{\theta})}{\partial \theta_1} \\ \vdots \\ \frac{\partial \ell(\underline{\theta})}{\partial \theta_b} \end{pmatrix} = \begin{pmatrix} \frac{\partial \log L(\underline{\theta}|\mathbf{Y})}{\partial \theta_1} \\ \vdots \\ \frac{\partial \log L(\underline{\theta}|\mathbf{Y})}{\partial \theta_b} \end{pmatrix}$$

This function is random because it depends on the data \mathbf{Y} .

Generally, the maximum likelihood estimator $\hat{\underline{\theta}}_{\text{MLE}}$ is the value of $\underline{\theta}$ where the maximum (over the parameter space Θ) of $L(\underline{\theta}|\mathbf{Y})$ is attained.

$$\hat{\underline{\theta}}_{\text{MLE}} = \underset{\underline{\theta} \in \Theta}{\operatorname{argmax}} L(\underline{\theta}|\mathbf{Y}) \iff L(\hat{\underline{\theta}}_{\text{MLE}}|\mathbf{Y}) \geq L(\underline{\theta}|\mathbf{Y}) \quad \forall \underline{\theta} \in \Theta$$

Under the assumption that the log-likelihood is continuously differentiable, then

$$S(\hat{\underline{\theta}}_{\text{MLE}}) = 0$$

But not always (!).

Example (Exponential threshold model): Suppose that Y_1, \dots, Y_n are iid from the exponential distribution with a threshold parameter μ ,

$$f(y; \mu) = \begin{cases} \exp\{-(y - \mu)\} & \mu \leq y < \infty \\ 0 & \text{otherwise,} \end{cases}$$

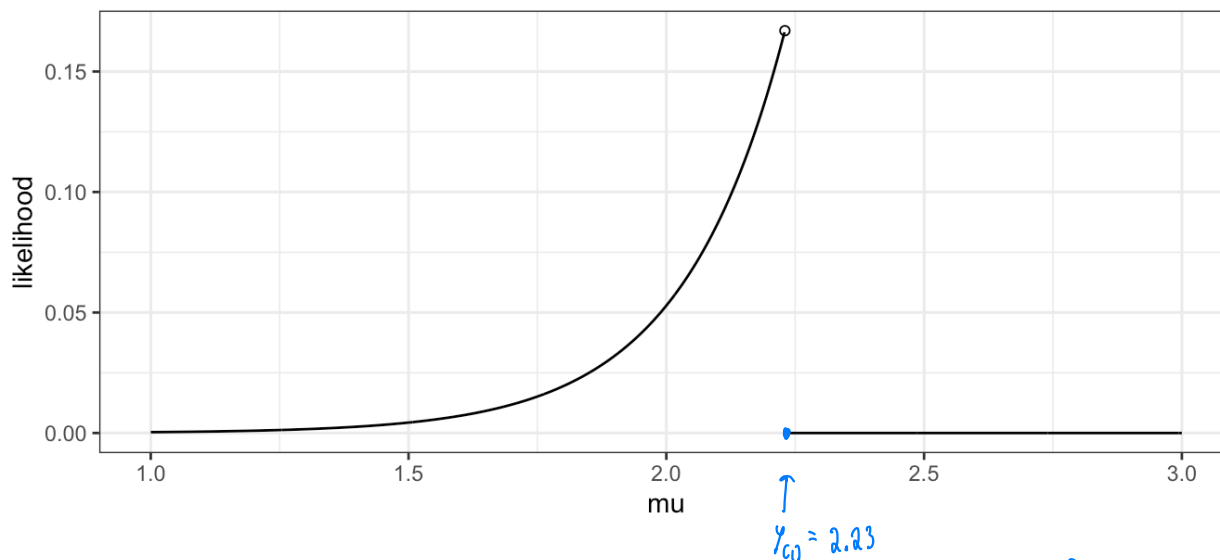
support depends on parameter.

for $-\infty < \mu < \infty$.

$$\begin{aligned} L(\mu | Y) &= \prod_{i=1}^n f(y_i; \mu) = \prod_{i=1}^n \exp(-(y_i - \mu)) \mathbb{I}(\mu < y_i) \\ &= \exp(-n\bar{y}) \exp(n\mu) \prod_{i=1}^n \mathbb{I}(\mu < y_i) \end{aligned}$$

\Rightarrow the likelihood $= 0$ for any value of $\mu \geq y_{(1)}$ $y_{(1)} = \min\{y_1, \dots, y_n\}$

Consider the artificial data set $\mathbf{y} = [2.47, 2.35, 2.23, 3.53, 2.36]$.



$\hat{\mu}_{MLE} = 2.23$ right? But $L(2.23 | Y) = 0 \Rightarrow L(\hat{\mu}_{MLE} | Y) \neq L(\mu | Y) \nmid \mu \in \mathbb{R}$ AND (*)

$S(\hat{\mu}_{MLE}) \neq 0$ because L not differentiable here.

If replace $\mu < y$ with $\mu \leq y$ in $f(y; \mu)$ then (*) will hold, but not the score equation.

To see this, Consider maximizing $L_h(\mu | Y) = \left(\frac{1}{2h}\right)^n \prod_{i=1}^n \{F_Y(y_i + h; \mu) - F_Y(y_i - h; \mu)\}$ for "small enough" value of h

Then maximize $\lim_{h \rightarrow 0^+} L_h(\mu | Y)$.

Rest of this section: assume support doesn't depend on parameter.

1.5.1 The Fisher Information Matrix

The Fisher information matrix $I(\theta)$ is defined as the $b \times b$ matrix where $\theta \in \mathbb{R}^b$

$$I_{ij}(\theta) = E \left[\left\{ \frac{\partial}{\partial \theta_i} \log f(\underline{y}_i; \underline{\theta}) \right\} \left\{ \frac{\partial}{\partial \theta_j} \log f(\underline{y}_i; \underline{\theta}) \right\} \right]$$

Note:

This is the "information" in one observation.

Is this random? No! It's an expectation.

In matrix form,

$$I(\theta) = E \left[\underbrace{\left(\frac{\partial}{\partial \theta^T} \log f(\underline{y}_i; \underline{\theta}) \right)}_{\text{column vector}} \underbrace{\left(\frac{\partial}{\partial \underline{\theta}} \log f(\underline{y}_i; \underline{\theta}) \right)}_{\text{row vector}} \right]$$

$$\text{let } \underline{s}(\underline{y}_i; \underline{\theta}) = \left\{ \frac{\partial}{\partial \underline{\theta}} \log f(\underline{y}_i; \underline{\theta}) \right\}^T \leftarrow \text{column vector}$$

$$\text{then } I(\theta) = E \left[\underline{s}(\underline{y}_i; \underline{\theta}) \underline{s}(\underline{y}_i; \underline{\theta})^T \right] \quad \text{"score contribution"}$$

again just depends
on 1 observation (not n)

Fisher information facts:

1. The Fisher information matrix is the variance of the score contribution.

Why? $E[s(y, \underline{\theta})] = 0$

fact ① from GLM section / Homework

Big Result

- ② If regularity conditions are met,

$$\underbrace{\sqrt{n}(\hat{\theta}_{\text{MLE}} - \underline{\theta})}_{\text{based on } n \text{ observations}} \xrightarrow{d} N_b(0, \underbrace{I(\underline{\theta})^{-1}}_{\text{inverse Fisher information}}).$$

defined on a single observation.

unbiased.

in practice

\Rightarrow If n is large,

$$\sqrt{n}(\hat{\theta}_{\text{MLE}} - \underline{\theta}) \overset{\text{"approximately distributed as"}}{\sim} N(\underline{0}, I(\underline{\theta})^{-1})$$

OR

$$\begin{aligned} \hat{\theta}_{\text{MLE}} - \underline{\theta} &\overset{\sim}{\sim} N(0, \frac{1}{n} I(\underline{\theta})^{-1}) \\ \hat{\theta}_{\text{MLE}} &\overset{\sim}{\sim} N(\underline{\theta}, \underbrace{\frac{1}{n} I(\underline{\theta})^{-1}}_{\{n I(\underline{\theta})\}^{-1}}) \end{aligned}$$

We will prove this result for $b=1$ next week.

3. If $b = 1$, then any unbiased estimator must have variance greater than or equal to $\{nI(\theta)\}^{-1}$

↑ Cramer-Rao lower bound.

If $b \geq 1$: If Σ is the asymptotic covariance matrix of any other consistent estimator then $\Sigma - I(\theta)^{-1}$ is positive definite.

4. The information matrix is related to the curvature of the log-likelihood contribution.

$$I(\theta) = E \left[\left\{ \frac{\partial}{\partial \theta} \log f(Y_i; \theta) \right\} \left\{ \frac{\partial}{\partial \theta} \log f(Y_i; \theta) \right\}^T \right]$$

$$= E \left[\left(\frac{\partial}{\partial \theta} \log f(Y_i; \theta) \right)^2 \right]$$

$$= E \left[- \frac{\partial^2}{\partial \theta^2} \log f(Y_i; \theta) \right]$$

$$= E \left[- \frac{\partial}{\partial \theta} s(Y_i; \theta) \right] \quad (\text{written a different way})$$

assuming l is twice differentiable and uses fact (2) from GLM section/homework.

Information is related to Hessian.

Comments:

(2) + (4) often combined to produce asymptotic confidence intervals for MLE's.

$$\text{Equivalence of } E \left[\left\{ \frac{\partial}{\partial \theta} \log f(Y_i; \theta) \right\} \left\{ \frac{\partial}{\partial \theta} \log f(Y_i; \theta) \right\}^T \right] \text{ and } E \left[- \frac{\partial^2}{\partial \theta^2} \log f(Y_i; \theta) \right]$$

relies on $Y_1, \dots, Y_n \stackrel{iid}{\sim} f(y; \theta)$

When this is not true, equality will not hold and (4) doesn't hold. (more later in the course)