

## 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 often restrict attention to likelihoods that are continuously differentiable wrt  $\theta$ .

In this case,

Recall the score function

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

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

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

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

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

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

But not always (?!).

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

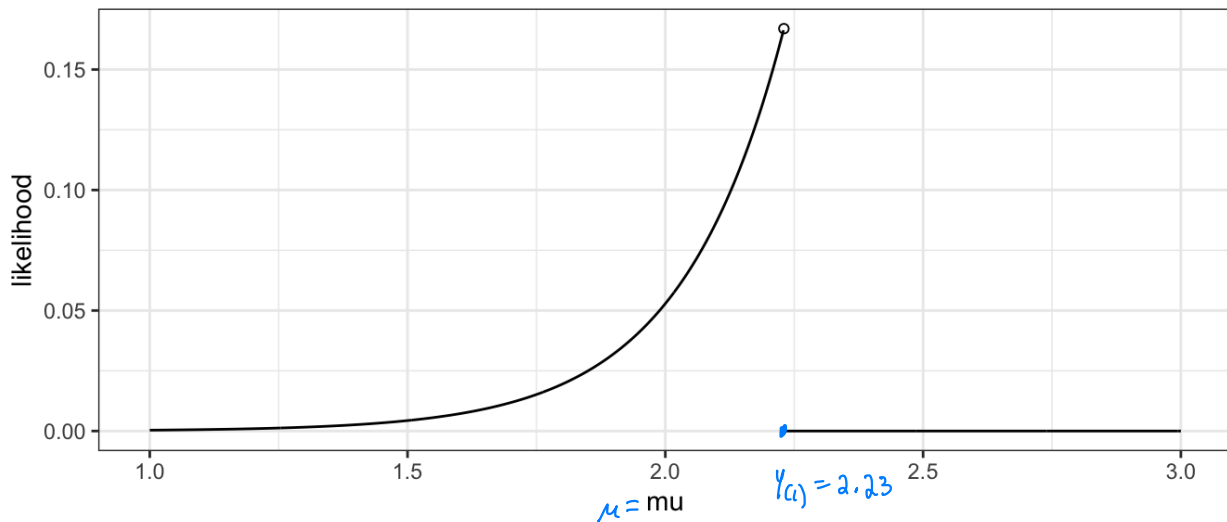
$$f(y; \mu) = \begin{cases} \exp\{-(y - \mu)\} & \mu \overset{\text{strict.}}{<} y < \infty \\ 0 & \text{otherwise,} \end{cases}$$

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

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

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

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

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

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

then maximize  $\lim_{h \rightarrow 0^+} L_h(\mu | \mathcal{Y}) \Rightarrow \hat{\mu}_{MLE} = Y_{(1)}$ .

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

### 1.5.1 The Fisher Information Matrix

$y_i \sim f(y; \theta)$

The Fisher information matrix  $I(\theta)$  is defined as the  $b \times b$  matrix where

dimension of  $\theta$

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

Is this random?  
No! it's an expectation!

Notice: this is the "information" in one observation. (note  $y_i$ ).

In matrix form,

$$I(\theta) = E \left[ \underbrace{\left( \frac{\partial}{\partial \theta} \log f(y_i; \theta) \right)}_{\text{column vector.}} \underbrace{\left( \frac{\partial}{\partial \theta} \log f(y_i; \theta) \right)}_{\text{row vector}} \right]$$

Let  $s(y; \theta) = \left\{ \frac{\partial}{\partial \theta} \log f(y; \theta) \right\}^T \leftarrow$  column vector.  
 $\uparrow$  score contribution.

Then  $I(\theta) = E \left[ s(y_i; \theta) s(y_i; \theta)^T \right].$

$\uparrow$   
Again this depends on 1 observation (not  $n$  of them).

Fisher information facts:

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

Why?  $E[S(Y_i, \theta)] = 0$

Fact ① from GLM section.

Big  
Result

- ② If regularity conditions are met,

$$\sqrt{n}(\hat{\theta}_{\text{MLE}} - \theta) \xrightarrow{d} N_b(\underline{0}, I(\theta)^{-1}).$$

based on  $n$   
observations.

↑ unbiased  
↖ inverse Fisher information.

defined wrt a single observation.

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

or,  $\hat{\theta}_{\text{MLE}} - \theta \overset{\sim}{\sim} N(\underline{0}, \frac{1}{n} I(\theta)^{-1})$

$\hat{\theta}_{\text{MLE}} \overset{\sim}{\sim} N(\theta, \{nI(\theta)\}^{-1})$  (\*)

We will prove this result for  $b=1$ . (later).

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

↪ Cramer-Rao lower bound.

If  $b \neq 1$ : If  $\Sigma$  is the asymptotic cov 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.  
Hessian

$$\begin{aligned}
 I(\theta) &= E \left[ \left( \frac{\partial}{\partial \theta} \log f(y_i; \theta) \right) \left( \frac{\partial}{\partial \theta} \log f(y_i; \theta) \right)^T \right] \\
 &= E \left[ - \frac{\partial^2}{\partial \theta \partial \theta^T} \log f(y_i; \theta) \right] \leftarrow \text{assuming } \log f \text{ is twice differentiable and} \\
 &\quad \text{using fact 2 from GLM section.} \\
 &= E \left[ - \frac{\partial}{\partial \theta} s(y_i; \theta) \right] \text{ (writing another way).}
 \end{aligned}$$

### 1.5.2 Observed Information

The information matrix is not random, but it is also not observable from the data.

Let  $Y_1, \dots, Y_n$  be iid with density  $f_Y(y_i; \boldsymbol{\theta})$ . The log likelihood is defined as

taking two derivatives and dividing by  $n$  results in

**Definition:** The matrix  $n\bar{I}(Y; \hat{\boldsymbol{\theta}}_{\text{MLE}})$  is called the sample information matrix, or the *observed information matrix*.

Why use  $I(\boldsymbol{\theta}) = \mathbf{E} \left[ -\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \log f(Y_1; \boldsymbol{\theta}) \right]$  as the basis for an estimator, rather than  $I(\boldsymbol{\theta}) = \mathbf{E} \left[ \left\{ \frac{\partial}{\partial \boldsymbol{\theta}^\top} \log f(Y_1; \boldsymbol{\theta}) \right\} \left\{ \frac{\partial}{\partial \boldsymbol{\theta}} \log f(Y_1; \boldsymbol{\theta}) \right\} \right]$ ?

Now let's prove the asymptotic normality of the MLE (in the scalar case).