# 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

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

*This function is random because it depends on the data $Y$.*

*?!*

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

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

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 $Y_1, \ldots, Y_n$ are iid from the exponential distribution with a threshold parameter $\mu$,
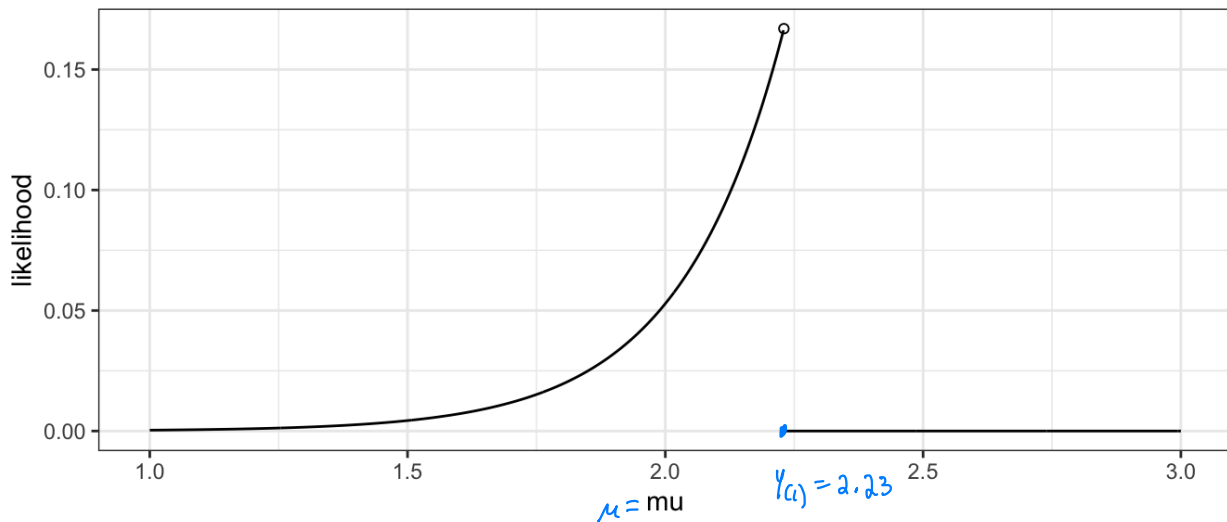
$$f(y; \mu) = \begin{cases} \exp\{-(y - \mu)\} & \mu \overset{\text{strict.}}{<} y < \infty \\ 0 & \text{otherwise,} \end{cases}$$

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

$$L(\mu \mid Y) = \prod_{i=1}^{n} f(Y_i; \mu) = \prod_{i=1}^{n} \exp\left(-(Y_i - \mu)\right) \mathbb{I}(\mu < Y_i).$$

$$= \exp(-n\bar{Y}) \exp(n\mu) \underbrace{\prod_{i=1}^{n} \mathbb{I}(\mu < Y_i)}_{\mathbb{I}(\mu < Y_{(1)})}.$$

$\Rightarrow$ the likelihood $= 0$ for any value of $\mu \geq Y_{(1)}$, $Y_{(1)} = \min\{Y_1, \ldots, Y_n\}$.

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



$\mu = $ mu          $Y_{(1)} = 2.23$

$\hat{\mu}_{MLE} = 2.23$, right?   $L(2.23 \mid Y) = 0$   $\Rightarrow (*) L(\hat{\mu}_{MLE} \mid Y) \not\geq L(\mu \mid Y) \; \forall \mu \in \mathbb{R}$.   AND

yes.

$S(\hat{\mu}_{MLE}) \neq 0$ because $\ell(\hat{\mu}_{MLE})$ 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 $L_h(\mu \mid Y) = \left(\frac{1}{2h}\right)^n \prod_{i=1}^{n} \{F_Y(Y_i + h; \mu) - F_Y(Y_i - h; \mu)\}$ for "small enough" value of $h$.

Then maximize $\lim_{h \to 0^+} L_h(\mu \mid Y)$.   $\Rightarrow \hat{\mu}_{MLE} = Y_{(1)}$.

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

## 1.5.1 The Fisher Information Matrix

$Y_i \overset{iid}{\sim} f(y; \theta)$

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

*dimension of $\theta$*

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

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

*Notice: this is the "information" in one observation. (note $Y_1$).*

In matrix form,

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

$\parallel$

Let $s(y; \underline{\theta}) = \left\{\frac{\partial}{\partial \underline{\theta}} \log f(y; \underline{\theta})\right\}^T$ ← *column vector.*

↳ *score contribution.*

Then $I(\underline{\theta}) = E\left[s(Y_i; \underline{\theta}) \, s(Y_i; \underline{\theta})^T\right]$.

*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? $E\left[S(Y_i, \theta)\right] = 0$

Fact ① from GLM section.

Big Result

2. If regularity conditions are met,

defined wrt a single observation.

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_{\text{MLE}} - \boldsymbol{\theta}) \xrightarrow{d} \text{N}_b(\mathbf{0}, I(\boldsymbol{\theta})^{-1}).$$

based on $n$ observations.

unbiased ~ inverse Fisher information.

"approximately distributed"

$\Rightarrow$ if $n$ is large $\quad \sqrt{n}\left(\hat{\underline{\theta}}_{MLE} - \underline{\theta}\right) \overset{\cdot}{\sim} N\left(\underline{0}, I(\theta)^{-1}\right)$

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

$\hat{\theta}_{MLE} \overset{\cdot}{\sim} N\left(\underline{\theta}, \{nI(\theta)\}^{-1}\right) \quad (*).$

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 $\{nI(\boldsymbol{\theta})\}^{-1}$

    ↰ Cramer-Rao lower bound.

If $b > 1$: If $\Sigma$ is the asymptotic cov matrix of any other consistent estimator, then

$$\Sigma - I(\theta)^{-1} \text{ is positive definite.}$$

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

    Hessian

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

$$= E\left[-\frac{\partial^2}{\partial\theta\,\partial\underline{\theta}^T} \log f(Y_i;\underline{\theta})\right] \quad \leftarrow \text{assuming } \ell \text{ is twice differentiable and using fact ② from GLM section.}$$

$$= E\left[-\frac{\partial}{\partial\theta} s(Y_i;\underline{\theta})\right] \text{ (writing another way).}$$

## 1.5.2 Observed Information

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

You need knowledge of the distribution to calculate it

Would be great to use $I(\hat{\theta}_{MLE}) = E\left\{ -\frac{\partial^2}{\partial\theta\partial\theta^T} \log f(Y_i; \underline{\theta})\Big|_{\underline{\theta}=\hat{\theta}_{MLE}} \right\}$

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

$$\log L(\underline{t}|\underline{y}) = \sum_{i=1}^{n} \log f_y(y_i; \underline{\theta})$$

taking two derivatives and dividing by $n$ results in

define $\quad \overline{I}(\underline{y}, \underline{\theta}) = \frac{1}{n}\sum_{i=1}^{n}\left\{ -\frac{\partial^2}{\partial\theta\partial\theta^T} \log f_y(y_i; \underline{\theta}) \right\}$

↖ average curvature contribution.

If $\quad I(\underline{\theta}) = E\left\{ -\frac{\partial^2}{\partial\theta\partial\theta^T} \log f(y_i; \underline{\theta}) \right\}$ then $\overline{I}(\underline{y}, \underline{\theta})$ would be an obvious estimator

$\times$ if we know $\underline{\theta}^*$

$\Rightarrow \overline{I}(\underline{y}, \hat{\theta}_{MLE})$ seems like a natural estimator for $I(\underline{\theta})$.

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

doesn't depend on sample size

Note: $I(\underline{\theta})$ is the expected curvature of the log-likelihood surface from <u>one</u> observation

The observed information $n\bar{I}(Y; \hat{\theta}_{\text{MLE}})$ is from a sample of size $n$ and does depend on sample size.

Recall $\hat{\theta}_{\text{MLE}} \sim N(\underline{\theta}, \{n I(\theta)\}^{-1})$ (*)

To get approximate variance of $\hat{\theta}_{\text{MLE}}$ for sample of size $n$, we need the matrix to depend on $n$.

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

The Hessian (curvature) @ $\hat{\theta}_{\text{MLE}}$ is readily available from optimization methods $\Rightarrow$

$n\bar{I}(Y, \hat{\theta}_{\text{MLE}})$ can be computed easily.

Alternatively could use $\bar{I}^*(Y, \theta) = \frac{1}{n}\sum_{i=1}^{n}\left[\left\{\frac{\partial}{\partial\underline{\theta}^\top}\log f(Y_i;\underline{\theta})\right\}\left\{\frac{\partial}{\partial\theta}\log f(Y_i;\theta)\right\}\right]$

$\left(\text{because } \text{E}\left[\bar{I}^*(Y,\theta)\right] = I(\theta) \text{ also}\right)$.

This is not typically used unless specification of $f$ is less clear (model misspecification).

$\left(\bar{I}(Y,\theta) \text{ is more efficient}\right)$

We will see this again later.

Now let's prove the asymptotic normality of the MLE (in the scalar case). $b=1$.

Useful facts: For $X_1, \dots, X_n$ iid $\operatorname{Var} X_1 = \sigma^2 < \infty$,

WLLN: $\overline{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i \xrightarrow{P} E[X_1].$

CLT: $\sqrt{n} \left( \overline{X}_n - EX_1 \right) = \sqrt{n} \left( \frac{1}{n} \sum_{i=1}^{n} X_i - EX_1 \right) \xrightarrow{d} N(0, \sigma^2).$

Let $Y_i \overset{iid}{\sim} f_Y(y; \theta)$ and $\hat{\theta}_{MLE}$ is such that $\frac{d}{d\theta} \ell(\theta) \big|_{\theta = \hat{\theta}_{MLE}} = 0.$
$S(\hat{\theta}_{MLE}).$

Let $S(\theta) = \frac{d}{d\theta} \ell(\theta) = \sum_{i=1}^{n} \frac{d}{d\theta} \log f(Y_i; \theta)$

$= \sum_{i=1}^{n} s(Y_i, \theta)$ where $s(Y_i, \theta) = \frac{d}{d\theta} \log f(Y_i; \theta)$

We know $E[s(Y_1, \theta)] = 0$ and $\operatorname{Var}[s(Y_1, \theta)] = I(\theta)$ and $\{s(Y_i, \theta)\}_{i=1}^{n}$ are iid r.v.'s.

$\Rightarrow \sqrt{n} \left( \frac{1}{n} S(\theta) - 0 \right) \xrightarrow{d} N(0, I(\theta))$ by CLT

$\iff (n I(\theta))^{-1/2} S(\theta) \xrightarrow{d} Z, \quad Z \sim N(0, 1) \quad (*)$

Secondly, let $J(\theta) = - \sum_{i=1}^{n} \frac{d^2 \log f_Y(Y_i; \theta)}{d\theta} = - \sum_{i=1}^{n} \frac{d}{d\theta} s(Y_i, \theta),$ Then $E\left[ -\frac{d}{d\theta} s(Y_1, \theta) \right] = I(\theta).$

$\underbrace{\qquad\qquad}_{\text{sum of iid R.V.'s}}$

$\Rightarrow \frac{1}{n} J(\theta) \xrightarrow{P} I(\theta)$ by WLLN $\iff n J^{-1}(\theta) \xrightarrow{P} I(\theta)^{-1} \quad (**)$

So far we have considered the true value $\theta$. Let $\ell(\theta)$ be sufficiently smooth to allow for Taylor Expansion.

$\overset{\text{assumption}}{\downarrow}$
$0 = S(\hat{\theta}_{MLE}) \approx S(\theta) + \frac{d S(\theta)}{d\theta} \left( \hat{\theta}_{MLE} - \theta \right) \iff \hat{\theta}_{MLE} - \theta \approx - \frac{1}{\frac{dS(\theta)}{d\theta}} \cdot S(\theta).$

$= \underline{J(\theta)^{-1} S(\theta).}$

The thing we want $\xrightarrow{d} N(0,1).$

Thus, $\overline{\sqrt{n} I(\theta)^{1/2} \left( \hat{\theta}_{MLE} - \theta \right)} \approx \sqrt{n} I(\theta)^{1/2} \underline{J(\theta)^{-1} S(\theta)}$

$= \underbrace{\{n I(\theta)\}^{1/2} J(\theta)^{-1} \{n I(\theta)\}^{1/2}}_{I(\theta) \, n \, J(\theta)^{-1}} \underbrace{\{n I(\theta)\}^{-1/2} S(\theta).}_{\xrightarrow{d} Z \,\, (*).}$

$\xrightarrow{P} I(\theta)^{-1} \quad (**),$

$\xrightarrow{d} N(0, 1)$ by Slutsky's Thm //

Note the argument to replace $I(\theta)$ by $I(\hat{\theta}_{MLE})$ in the asymptotic result is justified by convergence in probability.

The argument is generalized to $\underline{\theta}$ by interpreting the score as a $b \times 1$ vector, $I(\underline{\theta})$ as a $b \times b$ matrix, $Z \sim N_b(\underline{0}, I_b).$