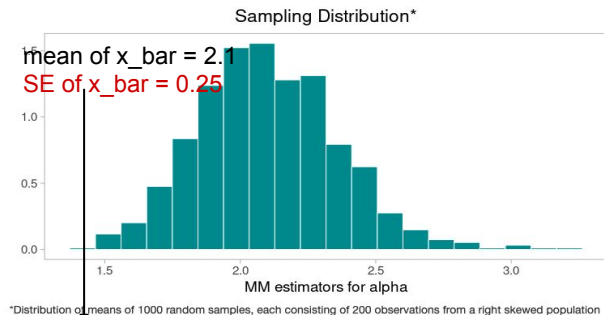


More Examples on MLE

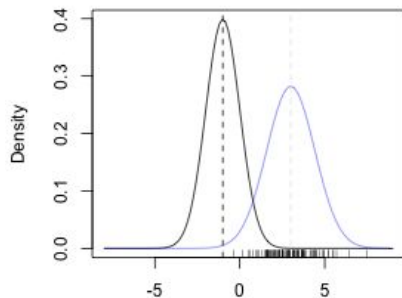
8.5 of Rice - Maximum Likelihood Estimators

06/29/2021

In the previous lecture,



Standard error of the estimate



Which one has^x more chance to generate such a sample?

- The Bootstrap simulation of the **sample distribution** of the MM estimator (Page 10 & 11):

- Mimicking the population distribution $f(x|\theta)$ by $f(x|\hat{\theta}_{MM})$
- Mimicking the theoretical SE $\sqrt{h(\theta, n)}$ by $\sqrt{h(\hat{\theta}_{MM}, n)}$.

This procedure can be generalized to MLE.

Coding practices in HW2 and Lab 3.

- Consistency of MM estimators:

$$\hat{\theta}_{MM} \xrightarrow{p} \theta, \text{ as } n \rightarrow \infty.$$

Compare with unbiasedness.

- Introduced maximum likelihood estimators (MLE):

- MM estimators can give unrealistic estimate when n is small.
- Maximize the likelihood over a meaning set of θ .
- Derived MLE for $N(\mu, \sigma^2)$ and $\text{Poisson}(\lambda)$.

Maximum likelihood estimators

Example 1. Let X_1, \dots, X_n be i.i.d Gamma(α, β). Find the MLE for α and β .

Under the i.i.d assumption,

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in \Theta} \prod_{i=1}^n f(X_i | \theta).$$

Solution. The probability density function is $f(x | \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}$.

$$L(\alpha, \beta) = \prod_{i=1}^n f(X_i | \alpha, \beta) = \frac{\beta^{n\alpha}}{\Gamma(\alpha)^n} \left(\prod_{i=1}^n X_i \right)^{\alpha-1} e^{-\beta \sum_{i=1}^n X_i}$$

$$\ell(\alpha, \beta) = \log L(\alpha, \beta) = n\alpha \log \beta - n \log \Gamma(\alpha) + (\alpha-1) \sum_{i=1}^n \log X_i - \beta \sum_{i=1}^n X_i$$

$$\begin{cases} \frac{\partial \ell}{\partial \alpha} = n \log \beta - n \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} + \sum_{i=1}^n \log X_i = 0 \\ \frac{\partial \ell}{\partial \beta} = n \frac{\alpha}{\beta} - \sum_{i=1}^n X_i = 0 \end{cases}$$

$$\beta = \frac{\alpha}{\bar{X}_n}$$

need to solve $\Rightarrow n \log \alpha - n \log \bar{X}_n - n \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} + \sum_{i=1}^n \log X_i = 0$

$$\hat{\alpha}_{MM} = 2.104$$

$$\hat{\beta}_{MM} = 2.298$$

Maximum likelihood estimators - think more about where to search $\downarrow \downarrow \begin{matrix} 0,1 \\ \end{matrix}$
Binomial (k, p)

\hat{k}_{ML} \hat{p}_{ML}

Definition. For the i.i.d samples X_1, \dots, X_n , we vary the value of θ in a meaningful set to evaluate its likelihood

$$L(\theta) = \prod_{i=1}^n f(X_i | \theta),$$

and $L(\theta)$ is called the likelihood function. The **maximum likelihood estimator** of θ is the particular value that maximizes the likelihood.

$\{x : f(x|\theta) > 0\}$ the support of density
 Maximum likelihood estimators function

$$X_{(1)} = \min \{X_1, \dots, X_n\}$$

$$X_{(n)} = \max \{X_1, \dots, X_n\}$$

Example 2. Let X_1, \dots, X_n be i.i.d $U(-\theta, \theta)$. Find the MM estimator and MLE for θ .

Solution: The density for $U(-\theta, \theta)$ is:

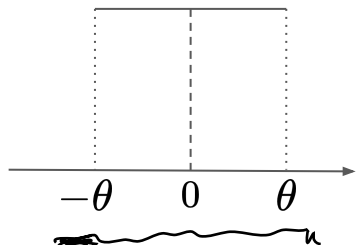
$$f(x|\theta) = \frac{1}{2\theta}, \quad x \in (-\theta, \theta)$$

$$\mu = EX = \int_{-\theta}^{\theta} \frac{x}{2\theta} dx = 0.$$

$$\mu_2 = EX^2 = \int_{-\theta}^{\theta} \frac{x^2}{2\theta} dx = \frac{\theta^3}{3}.$$

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n X_i^2 = \frac{\hat{\theta}_{MM}^3}{3}$$

$$\Rightarrow \hat{\theta}_{MM} = \left(\frac{3}{n} \sum_{i=1}^n X_i^2 \right)^{1/3}$$



For MLE, let's look into the set where θ is likely to generate a sample of X_1, \dots, X_n .

$$-\theta \leq X_i \leq \theta, \quad i=1, \dots, n$$

$$\Leftrightarrow -\theta \leq X_{(1)} < X_{(n)} \leq \theta$$

$$\Leftrightarrow \theta \geq -X_{(1)}, \quad \theta \geq X_{(n)}$$

$$\Leftrightarrow \theta \geq \max \{-X_{(1)}, X_{(n)}\}$$

$$L(\theta) = \prod_{i=1}^n f(X_i|\theta) = \left(\frac{1}{2\theta} \right)^n$$

$$\hat{\theta}_{MLE} = \max \{-X_{(1)}, X_{(n)}\}$$

Maximum likelihood estimators

$$\hat{\theta}_{MLE} = \max \{ -X_{(1)}, X_{(n)} \}$$



Example 2 cont'd. Let X_1, \dots, X_n be i.i.d $U(-\theta, \theta)$. Find the MLE for θ .

Can you work out the theoretical sampling distribution of $\hat{\theta}_{MLE}$? (3.7 of Rice)

$V = X_{(1)}$ and $U = X_{(n)}$ have the joint density

$$\begin{aligned} f(u, v) &= n(n-1) f(v) f(u) [F(u) - F(v)]^{n-2}, \quad u \geq v \\ &= n(n-1) \frac{1}{2\theta} \frac{1}{2\theta} \left[\frac{u-v}{2\theta} \right]^{n-2} \\ &= \frac{n(n-1)}{8\theta^3} (u-v)^{n-2} \end{aligned}$$

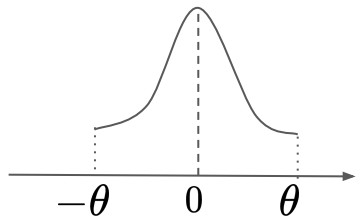
$$\begin{aligned} P(\hat{\theta}_{MLE} \leq t) &= P(\max \{ -X_{(1)}, X_{(n)} \} \leq t) = P(-X_{(1)} \leq t, X_{(n)} \leq t) \\ &= P(\underline{X_{(1)} \geq -t, X_{(n)} \leq t}) = \int_{-t}^{\theta} \int_{-t}^t f(u, v) du dv \end{aligned}$$

Maximum likelihood estimators

Example 3. Let X_1, \dots, X_n be i.i.d from a population with density

$$f(x|\theta) = \begin{cases} \frac{3x^2}{2\theta^3} & \text{if } -\theta \leq x \leq \theta, \\ 0, & \text{otherwise.} \end{cases}$$

$$\begin{aligned} f(u, v) &= n(n-1) f(u) f(v) (F(u) - F(v))^{n-2} \\ &= n(n-1) \frac{3u^2}{2\theta^3} \frac{3v^2}{2\theta^3} \int_v^u \frac{3t^2}{2\theta^3} dt \end{aligned}$$



Solution: $-\theta \leq X_1, \dots, X_{(n)} \leq \theta$

$$\Leftrightarrow \theta \geq \max \{ -X_{(1)}, X_{(n)} \}$$

$$L(\theta) = \prod_{i=1}^n f(x_i|\theta) = \prod_{i=1}^n \frac{3x_i^2}{2\theta^3} = \frac{3^n \prod_{i=1}^n x_i^2}{2^n \theta^{3n}}$$

$$\hat{\theta}_{MLE} = \max \{ -X_{(1)}, X_{(n)} \}$$

Large Sample Theories for MLE

8.5.2 of Rice

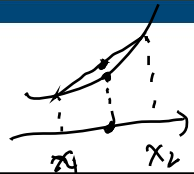
06/29/2021

How good are the MLE?

Jensen's inequality:

ϕ convex

$$\phi(x_1) + \phi(x_2) \geq 2\phi\left(\frac{x_1 + x_2}{2}\right)$$



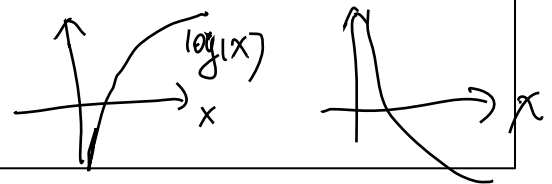
ϕ convex

$$\Rightarrow E\phi(X) \geq \phi(E(X))$$

Lemma A. Suppose $X \sim f(x|\theta_0)$. For any other θ , the average difference between $\log f(X|\theta_0)$ and $\log f(X|\theta)$ is called the Kullback-Leibler divergence. Moreover,

$$KL(\theta_0, \theta) = E_{\theta_0} \left(\log \frac{f(X|\theta_0)}{f(X|\theta)} \right) \geq 0.$$

The KL divergence equals to zero only when $f(x|\theta) \equiv f(x|\theta_0)$.



A measure of how one probability distribution is different from a second, reference probability distribution, or relative entropy. $-\log(x)$

Proof. Since $-\log(x)$ is a convex function of x , we know by Jensen's inequality that

$$\begin{aligned} -E_{\theta_0} \left[\log \frac{f(X|\theta_0)}{f(X|\theta)} \right] &\geq -\log \left[E_{\theta_0} \frac{f(X|\theta_0)}{f(X|\theta)} \right] \\ &= -\log \left[\int f(x|\theta_0) \frac{f(x|\theta_0)}{f(x|\theta)} dx \right] \\ &= -\log \left[\int f(x|\theta) dx \right] \\ &= -\log 1 = 0. \end{aligned}$$

$f(x|\theta)$ is a density

$N(\mu, \sigma^2)$
 $f(x|\mu, \sigma^2)$
 $e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
 $\mu_1 \neq \mu_2$ will give different densities.

Identifiable if different θ 's produce different densities.

How good are the MLE?

Theorem B. Under the i.i.d and a few other assumptions, $\hat{\theta}_{MLE}$ is a consistent estimator of θ :

$$\hat{\theta}_{MLE} \xrightarrow{p} \theta, \text{ as } n \rightarrow \infty.$$

Proof*. Recall that the log-likelihood function is $l(\theta) = \log \left(\prod_{i=1}^n f(X_i | \theta) \right) = \sum_{i=1}^n \log f(X_i | \theta)$.

1. $f(x | \theta)$ is identifiable; \leftarrow
2. $l(\theta)$ is differentiable;
3. $\{x : f(x | \theta) > 0\}$ does not depend on θ .

$\frac{1}{n} \sum_{i=1}^n \log f(X_i | \theta)$ Law of large numbers \Leftarrow if X_1, \dots, X_n are i.i.d then $f(X_1 | \theta), \dots, f(X_n | \theta)$ are also i.i.d

$\xrightarrow{P} E_{\theta} \log f(X | \theta)$

$$\left(\frac{1}{n} l(\theta) - \frac{1}{n} l(\theta_0) \right) \approx E_{\theta} \log f(X | \theta) - E_{\theta_0} \log f(X | \theta_0)$$

$$= E_{\theta_0} \log \frac{f(X | \theta)}{f(X | \theta_0)} = - \text{KL}(\theta_0, \theta) \leq 0.$$

$$\frac{\max_{\theta} l(\theta)}{n} - \frac{1}{n} l(\theta_0) \leq 0$$

$$\hat{\theta}_{MLE} \approx \theta_0$$

The equality is only attained when $f(x | \theta) \equiv f(x | \theta_0) \Rightarrow \theta = \theta_0$ identifiability

If $\hat{\theta}_n$ is MLE, then $l'(\hat{\theta}_n) = 0$.

How good are the MLE? - Stronger than consistency

Denote $l_i(\theta) = \log f(X_i | \theta)$, then $l(\theta) = \sum_{i=1}^n l_i(\theta)$.

Use Taylor expansion to examine $l'_i(\theta)$.

$$l'_i(\hat{\theta}_n) - l'_i(\theta_0) = l''_i(\theta_0) (\hat{\theta}_n - \theta_0) + e_i$$

infinitely small

Sum over i

$$\Rightarrow \sum_{i=1}^n l'_i(\hat{\theta}_n) - \sum_{i=1}^n l'_i(\theta_0) = \sum_{i=1}^n l''_i(\theta_0) \cdot (\hat{\theta}_n - \theta_0) + \sum_{i=1}^n e_i$$

E

$$\Rightarrow l'(\hat{\theta}_n) - l'(\theta_0) = l''(\theta_0) \cdot (\hat{\theta}_n - \theta_0) + E$$

$$\Rightarrow \hat{\theta}_n - \theta_0 = \frac{-(l'(\hat{\theta}_n) - l'(\theta_0) + E)}{-l''(\theta_0)}$$

$$\Rightarrow \sqrt{n}(\hat{\theta}_n - \theta_0) = \frac{\sqrt{n}l'(\theta_0) - \sqrt{n}l'(\hat{\theta}_n) + E_n}{-l''(\theta_0)} \approx \frac{\frac{1}{n}\sqrt{n}l'(\theta_0)}{-\frac{l''(\theta_0)}{n}}$$

Numerator

$$\frac{1}{\sqrt{n}}l'(\theta_0) = \sqrt{n} \times \frac{1}{n} \sum_{i=1}^n \left(\frac{\partial}{\partial \theta} \log f(X_i | \theta_0) \right)$$

$$\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(X_i | \theta_0) - E \frac{\partial}{\partial \theta} \log f(X | \theta_0) \right)$$

$= 0$

CLT

$$\rightarrow N\left(0, \text{var} \frac{\partial}{\partial \theta} \log f(X | \theta_0)\right)$$

$= I(\theta_0)$

Denominator

$$-\frac{1}{n}l''(\theta_0) = -\frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log f(X_i | \theta_0)$$

$$\rightarrow -E_{\theta_0} \frac{\partial^2}{\partial \theta^2} \log f(X_i | \theta_0)$$

$= I(\theta_0)$

How good are the MLE? - One parameter model

Definition. The Fisher information of θ is defined by

$$I(\theta) = E_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X|\theta) \right]^2.$$

Lemma C. Under appropriate conditions on $f(x|\theta)$, the Fisher information can also be written as

The exchangeability of
differentiation & integration

$$I(\theta) = -E_{\theta} \left[\frac{\partial^2}{\partial \theta^2} \log f(X|\theta) \right].$$

(*) We need to show that $E_{\theta} \left[\frac{\partial^2}{\partial \theta^2} \log f(X|\theta) \right] = 0$

Note $E_{\theta} \left[\frac{\frac{\partial^2}{\partial \theta^2} f(x|\theta)}{f(x|\theta)} \right] = \int \frac{\frac{\partial^2}{\partial \theta^2} f(x|\theta)}{f(x|\theta)} f(x|\theta) dx$

$= \int \frac{\partial^2}{\partial \theta^2} f(x|\theta) dx$

$= \frac{\partial^2}{\partial \theta^2} \int f(x|\theta) dx = 0$

Proof. We begin by examining the 1st / 2nd derivative of $\log f(x|\theta)$:

$$\frac{\partial}{\partial \theta} \log f(x|\theta) = \frac{\frac{\partial}{\partial \theta} f(x|\theta)}{f(x|\theta)}, \quad \frac{\partial^2}{\partial \theta^2} \log f(x|\theta) = \frac{\frac{\partial^2}{\partial \theta^2} f(x|\theta) \cdot f(x|\theta) - \left[\frac{\partial f(x|\theta)}{\partial \theta} \right]^2}{f^2(x|\theta)}$$

$$= \frac{\frac{\partial^2}{\partial \theta^2} f(x|\theta)}{f(x|\theta)} - \left[\frac{\partial}{\partial \theta} \log f(x|\theta) \right]^2$$

Therefore,

$$E_{\theta} \left[\frac{\partial^2}{\partial \theta^2} \log f(X|\theta) \right] = E_{\theta} \left[\frac{\frac{\partial^2}{\partial \theta^2} f(X|\theta)}{f(X|\theta)} \right] - E_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X|\theta) \right]^2$$

(*)

$= 0 - I(\theta)$

see above

Cont'd from Page 11:

Using the results from Page 11, we know

$$\sqrt{n} \left\{ \frac{1}{n} \ell'(\theta) - E_{\theta_0} \left[\frac{\partial}{\partial \theta} \log f(X|\theta_0) \right] \right\} \rightarrow N(0, \text{var}_{\theta_0} \left(\frac{\partial}{\partial \theta} \log f(X|\theta_0) \right)).$$

Note: $E_{\theta_0} \left[\frac{\partial}{\partial \theta} \log f(X|\theta_0) \right] = \int \frac{\frac{\partial}{\partial \theta} f(X|\theta_0)}{f(X|\theta_0)} f(X|\theta_0) dX$

Another
exchangeability
assumption

$$= \int \frac{\partial}{\partial \theta} f(X|\theta_0) dX = \frac{\partial}{\partial \theta} \int f(X|\theta_0) dX = 0.$$

$$\text{var}_{\theta_0} \left[\frac{\partial}{\partial \theta} \log f(X|\theta_0) \right] \stackrel{\text{var } X = EX^2 - (EX)^2}{=} E_{\theta_0} \left[\frac{\partial}{\partial \theta} \log f(X|\theta_0) \right]^2 - \left\{ E_{\theta_0} \left[\frac{\partial}{\partial \theta} \log f(X|\theta_0) \right] \right\}^2$$

$= I(\theta_0)$ by definition $= 0$

$$= I(\theta_0).$$

Therefore, $\frac{1}{\sqrt{n}} \ell'(\theta) \rightarrow N(0, I(\theta_0))$

Lemma C

For the denominator, $-\frac{1}{n} \sum_{i=1}^n \ell''(\theta_0) \xrightarrow{P} E_{\theta_0} \left[\frac{\partial^2}{\partial \theta^2} \log f(X|\theta_0) \right] = I(\theta_0)$

Asymptotic Normality of MLE - One parameter model



$$P(|\hat{X}_n - \theta| > \epsilon) \rightarrow 0$$

To sum up,

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \approx \frac{\frac{1}{\sqrt{n}}l'(\theta_0)}{-\frac{1}{n}l''(\theta_0)} \approx$$

in which

$$\frac{1}{\sqrt{n}}l'(\theta_0) \rightarrow N(0, I(\theta_0)),$$

and

$$-\frac{1}{n}l''(\theta_0) = -\frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log f(X_i | \theta_0) \xrightarrow{p} I(\theta_0).$$

$$\frac{N(0, I(\theta_0))}{I(\theta_0)}$$

$$N(0, \frac{1}{I(\theta_0)})$$

1. $f(x | \theta)$ is identifiable;
2. $l(\theta)$ is differentiable; *Exchangeability*
3. $\{x : f(x | \theta) > 0\}$ *Assumption*
does not depend on θ ; \downarrow
4. $\int \frac{\partial}{\partial \theta} f(x | \theta) dx = \frac{\partial}{\partial \theta} \int f(x | \theta) dx,$
 $\int \frac{\partial^2}{\partial \theta^2} f(x | \theta) dx = \frac{\partial^2}{\partial \theta^2} \int f(x | \theta) dx;$
This can be easily satisfied by sufficiently smooth densities.

Theorem D. Under the i.i.d and a few other assumptions, the MLE $\hat{\theta}_n$ has asymptotic normality:

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N\left(0, \frac{1}{I(\theta_0)}\right).$$

Rate of convergence: $n^{-1/2}$
 $SE(\hat{\theta}_n) = \sqrt{\frac{1}{n I(\theta_0)}} = (n I(\theta_0))^{-1/2}$

$$\hat{\theta}_n - \theta_0 \rightarrow N\left(0, \frac{1}{n I(\theta_0)}\right)$$

