

Chapter 6 - Likelihood Inference

Section 6.1. The likelihood function.

Data s and the model $\{P_\theta : \theta \in \Omega\}$.

Start with P_θ that is discrete, given by $f_\theta : \Omega \rightarrow \mathbb{R}^+$.

Consider $L(\cdot|s) : \Omega \rightarrow \mathbb{R}^+$ s.t. $L(\theta|s) = f_\theta(s)$ —

θ : Data are fixed, but θ varies. Here, $f_\theta(s)$ is the probability of obtaining s when the true value of the parameter is

θ . Belief ordering on Ω : we prefer θ_1 over θ_2 if

$f_{\theta_1}(s) > f_{\theta_2}(s)$. Indifferent between θ_1 and θ_2 if $f_{\theta_1}(s) = f_{\theta_2}(s)$.

Interpretation: $L(\theta|s)$ is the probability of s if θ is the true value.

It may be very small but it is the ordering that is important, not the absolute value of $L(\theta|s)$.

Ex 6.1.1. Let $S = \{1, 2, \dots\} \times \{P_\theta : \theta \in \Omega\}$.

P_1 is $\sim \text{Unif}\{1, \dots, 10^3\}$ & P_2 is $\text{Unif}\{1, \dots, 10^6\}$. Observes = 10. Then $L(1/10) = \frac{1}{10^3}$ and $L(2/10) = \frac{1}{10^6}$. $\theta = 1$ is a 1000 times

more likely under Model 1 than under Model 2.

Thus, we're interested only in Likelihood ratios:

$\frac{L(\theta_1|s)}{L(\theta_2|s)}$ for $\theta_1, \theta_2 \in \Omega$. Thus, for any $c > 0$,

the function $L^*(\cdot|s) = cL(\cdot|s)$ can also serve as the likelihood function. $L(\cdot|s)$ and $L^*(\cdot|s)$ are two equivalent likelihoods.

Ex. 6.1.2. A coin is tossed $n=10$ times, 324 heads are observed.

The appropriate model for the # of heads is $\text{Bin}(10, \theta)$, $\theta \in \Omega = \{0, 1\}$.

$\Rightarrow L(\theta|s) = \binom{10}{4} \theta^4 (1-\theta)^6$, $\theta = 0.4$ is the maximum with the value 0.2508.

There are different approaches to inference via the likelihood function. At one extreme is the likelihood principle: if two model and data combinations yield equivalent likelihood functions, then inferences about the unknown parameter must be the same. This is a very restrictive approach.

Ex. 6.1.3. Toss a coin (independently) until 4 heads are obtained; the number of tails observed is $s = 6$.
 $S \sim \text{NegBin}(4, \theta)$ and $L(\theta|s) = \binom{9}{6} \theta^4 (1-\theta)^6$. This is a positive multiple of $L(\theta|4)$ considered in the Ex. (6.1.2)
 Note that the data were obtained in the entirely different ways!! ... additional model features besides the likelihood function must be taken into account.

An example of the likelihood inference: let $C(s) = \{\theta : L(\theta|s) > c_0\}$ for some $c_0 > 0$ - a likelihood region. This may possibly contain the true value θ ; if some $\theta^* \notin C(s)$, then $L(\theta^*|s) < L(\theta|s)$ for $\forall \theta \in C(s)$ and so θ^* is not well supported by the data compared to a $\theta \in C(s)$.
 For now, we don't know how to choose C - leave it for later.

If P_θ is continuous (imagine $f_{\theta_1}(s) > f_{\theta_2}(s)$ and $s \in \mathbb{R}^I$).

Then, if $f_\theta(s)$ is continuous at every s , we have

$$P_{\theta_1}(V) = \int f_{\theta_1}(s) ds > P_{\theta_2}(V) = \int f_{\theta_2}(s) ds \text{ for every } V = (a, b)$$

that is small enough and contains s the data support θ_1 more than θ_2 . If, $s \in \mathbb{R}^n$, $n > I$ - similar interpretation.

Thus, $L(\theta|s) = f_\theta(s)$ again.

Ex. 6.1.4. Location Normal model.

$(x_1, \dots, x_n) \sim N(\theta, \sigma_0^2)$ where $\sigma_0^2 > 0$ is known and, say, $\theta \in \mathbb{R} = \mathbb{R}^I$.

$$\text{Then, } L(\theta|x_1, \dots, x_n) = \prod_{i=1}^n f_\theta(x_i) = (2\pi\sigma_0^2)^{-n/2} e^{-\frac{n}{2\sigma_0^2}(\bar{x}-\theta)^2} = \frac{1}{(2\pi\sigma_0^2)^{n/2}} e^{-\frac{n}{2\sigma_0^2}(\bar{x}-\theta)^2}$$

The simpler version is $L(\theta|x_1, \dots, x_n) = \ell = e^{-\frac{n}{2\sigma_0^2}(\bar{x}-\theta)^2}$

If $n=25$, $\hat{\theta}_0^2 = 1$, observe $\bar{x} = 3.3 \dots$ see Fig. (6.12)

The peak is at $\theta = \bar{x} = 3.3$. The likelihood interval

$$C(x) = \{ \theta : L(\theta | x_1, \dots, x_n) \geq 0.5 \} = [3.0645, 3.53578]$$

contains all those θ for which likelihood is at least 0.5.

Ex. 6.1.5. Multinomial models.

Suppose that $k=3$, $(\theta_1, \theta_2, \theta_3)$ are unknown.

$$\mathcal{R} = \{ (\theta_1, \theta_2, \theta_3) : \theta_i \geq 0, i=1, 2, 3, \sum \theta_i = 1 \}.$$

\mathcal{R} is effectively two-dimensional!

We observe $(s_1, s_2, s_3) \Rightarrow L(\theta_1, \theta_2, \theta_3 | s_1, s_2, s_3) =$
 $= \theta_1^{x_1} \theta_2^{x_2} \theta_3^{x_3}$ where x_i is the number of i 's in the sample. We can see that the likelihood based on

the observed counts (x_1, x_2, x_3) is given by

$$L(\theta_1, \theta_2, \theta_3 | x_1, x_2, x_3) = \theta_1^{x_1} \theta_2^{x_2} \theta_3^{x_3}$$

Section 6.1.1 Sufficient statistics.

Def. A function T defined on the sample space

S is a sufficient statistic for the model \mathcal{T} , whenever $T(s_1) = T(s_2)$, then

$$L(\cdot | s_1) = c(s_1, s_2) L(\cdot | s_2) \text{ for some } c(s_1, s_2) > 0.$$

Thus, if we observe $T=t$, we can pick any value $s \in T^{-1}(t) = \{s : T(s)=t\}$ and base the likelihood on it. This is, effectively, a dimension reduction technic.

Ex. Take $S = \{1, 2, 3, 4\}$, $\Omega = \{\alpha, \beta\}$

$\theta = \alpha$	$\frac{1}{2}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$;
$\theta = \beta$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$;

$L(\cdot | 2) = L(\cdot | 3) = L(\cdot | 4)$, so $\{2, 3, 4\}$ give the same likelihood ratios. Thus, $T: S \rightarrow \{0, 1\}$ s.t. $T(1) = 0$, $T(2) = T(3) = T(4) = 1$ is a sufficient statistic. Now, the sample space for T has only two elements.

Factorization Theorem If $f_\theta(s) = h(s) g_\theta(T(s))$, $g_\theta, h \geq 0$, then T is a sufficient statistic.

Proof. If $T(s_1) = T(s_2)$,

$$L(\cdot | s_1) = h(s_1) g_\theta(T(s_1)) = \frac{h(s_1) g_\theta(T(s_1))}{h(s_2) g_\theta(T(s_2))} h(s_2) g_\theta(T(s_2)) =$$

$$= \frac{h(s_1)}{h(s_2)} h(s_2) g_\theta(T(s_2)) = c(s_1, s_2) L(\cdot | s_2) \text{ because}$$

$$g_\theta(T(s_1)) = g_\theta(T(s_2)).$$

A sufficient statistic hinges on the choice of parameter – if the parameter is different, $T(\cdot)$ will also change!!

Def. 6.1.2 A sufficient statistic T for a model is a minimal sufficient statistic, whenever the value of $T(s)$ can be calculated once we know the likelihood $f \sim L(\cdot|s)$.

A minimal sufficient statistic gives the greatest reduction of the data in the sense that, if T is a ms. suff. and U is suff., then there $\exists f \sim h$ such that $T = h(U)$. Sufficiency always depends on the model.

Ex. Location Normal Model

Factorization Theorem $\rightarrow \bar{x}$ is a sufficient statistic.

Any $L(\cdot|s)$ is a positive multiple of $\exp\left(-\frac{n}{2\sigma^2}(\bar{x}-\theta)^2\right)$

Any such function of θ is completely specified by the maximum point, that is, $\theta = \bar{x} \Rightarrow \bar{x}$ is a complete sufficient statistic (because $\sigma^2 > 0$ is known!).

Ex. Location - Scale Normal model

$$x_1, \dots, x_n \sim N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0 \text{ are unknown.}$$

$$\theta = (\mu, \sigma^2) \in \Omega = \mathbb{R}^2 \times (0, \infty) \Rightarrow$$

$$\Rightarrow L(\theta | x_1, \dots, x_n) = (2\pi\sigma^2)^{-n/2} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^{n-1} (x_i - \mu)^2} =$$

$$= (2\pi\sigma^2)^{-n/2} e^{-\frac{1}{2\sigma^2} (\bar{x} - \mu)^2 - \frac{n-1}{2\sigma^2} s^2} \Rightarrow (\bar{x}, s^2) \text{ is a suff. stat.}$$

If σ^2 is fixed, any positive multiple of $L(\cdot | x_1, \dots, x_n)$ is maximized at $\mu = \bar{x}$. Fix μ at \bar{x} , we have $= \frac{n-1}{2\sigma^2} s^2$

$$L((\bar{x}, \sigma^2) | x_1, \dots, x_n) = (2\pi\sigma^2)^{-n/2} e^{-\frac{n-1}{2\sigma^2} s^2} - \max$$

doesn't change if taking a log; \Rightarrow

$$\frac{\partial \ln L((\bar{x}, \sigma^2) | x)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{n-1}{2\sigma^4} s^2 / = 0$$

gives $\hat{\sigma}^2 = \frac{n-1}{n} s^2$ - it is a 1-1 function of s^2 .

$\Rightarrow (\bar{x}, s^2)$ is a minimal sufficient statistic for this model.

Ex., Multinomial model. - 6-

The likelihood function $L(\theta_1, \theta_2, \theta_3 | S_1, S_2, S_3) = \theta_1^{x_1} \theta_2^{x_2} \theta_3^{x_3}$) x_i - number of i's in the sample.

It is maximized by taking $(\hat{\theta}_1, \hat{\theta}_2, \hat{\theta}_3) = \left(\frac{x_1}{n}, \frac{x_2}{n}, \frac{x_3}{n} \right)$ \rightarrow

\rightarrow given the likelihood, we can compute counts if n is assumed known $\rightarrow (x_1, x_2, x_3)$ is a maximal sufficient statistic.

6.2 MLE - estimation problem.

Need to find $\hat{\theta}(s)$ that maximizes $L(\theta|s)$, i.e.

$$(*) L(\hat{\theta}(s)|s) \geq L(\theta|s) \text{ for } \forall \theta \in \Omega$$

$\hat{\theta}: S \rightarrow \Omega$ satisfying is an MLE; the value $\hat{\theta}(s)$ is called the maximum likelihood estimate

Ex. Let $S = \{S_1, S_2, S_3\}$; $\Omega = \{f_{1,2}\}$,

$$f_1(s) \quad 0.3 \quad 0.4 \quad 0.3$$

$$f_2(s) \quad 0.1 \quad 0.7 \quad 0.2$$

Observe $s=1$ (chips labelled 1, 2, 3 in given proportions in two bowls; how to select a bowl if a chip labelled

3 has been selected?)

$$\Rightarrow \hat{\theta}(1) = 1, \text{ because } 0.3 = L(1|1) > 0.1 = L(2|1).$$

$$\text{If we observed } s=2, \text{ then } \hat{\theta}(2) = 2; \text{ if } s=3, \hat{\theta}(3) = 1.$$

Remarks (i) MLE is not unique; if $f_2(1) = 0$, $f_2(2) = 0.7$, $f_2(3) = 0.3$, then, in addition to $\hat{\theta}(3) = 1$ we also have $\hat{\theta}(3) = 2$.

(ii) MLE is invariant ... Instead of $\theta \in \Omega$ use

$\gamma \in \Gamma = f(\theta): \theta \in \Omega \}$. E.g. take $f(1) = a, f(2) = b$, so that $\gamma = \{a, b\}$. The model is $\{g_\gamma: \gamma \in \Gamma\}$ where $g_\gamma = f_\theta$ for the unique value θ s.t. $f(\theta) = \gamma$. There is a new parameter $\gamma \in \Gamma$ - but there are labels only! Thus .

Th. If $\hat{\theta}(s)$ is an MLE when $\theta \in \Omega$, ψ is a 1-1 function defined on Ω , then $\hat{\psi}(s) = \psi(\hat{\theta}(s))$ is an MLE in the new parameterization.

Proof. If we select $L^*(\psi(s)) = g_\psi(s)$, and the original one is $L(\theta/s) = f_\theta(s)$, then

$$L^*(\hat{\psi}(s)/s) = g_{\hat{\psi}}(\hat{\theta}(s))(s) = f_{\hat{\theta}(s)}(s) = \varphi(\hat{\theta}(s)/s) \geq$$

$$\geq L(\theta/s) = L^*(\psi(\theta)/s) \text{ for } \forall \theta \in \Omega. \Rightarrow$$

$$\Rightarrow L^*(\hat{\psi}(s)/s) \geq L^*(\psi(s)) \quad \forall \psi \in \Psi \text{ q.e.d.}$$

6.2.1 Computation of the MLE

Definition. $L(\cdot/s) \Rightarrow l(\cdot/s)$ is $l(\cdot/s) = \ln L(\cdot/s)$

because \ln is an increasing function, this works.

Indeed, $L(\theta/s_1, s_n) = \prod_{i=1}^n f_\theta(s_i) \Rightarrow l(\theta/s_1, s_n) = \sum_{i=1}^n \ln f_\theta(s_i)$

Define $S(\theta/s) = \frac{\partial l(\theta/s)}{\partial \theta}$ - the score function.

To obtain the MLE have to solve $S(\theta/s) = 0$; to guarantee the solution is at least a ~~global~~ local max., check that $\frac{\partial S(\theta/s)}{\partial \theta}|_{\theta=\hat{\theta}(s)} =$

$$= \frac{\partial^2 l(\theta/s)}{\partial \theta^2}|_{\theta=\hat{\theta}(s)} < 0.$$

Ex Location normal model. Here $L(\theta/x_1, x_n) = e^{-\frac{n}{2\sigma_0^2}(\bar{x}-\theta)^2}$; check that $S(\theta/x_1, x_n) = \frac{n}{\sigma_0^2}(\bar{x}-\theta)|_{\theta=\hat{\theta}(x_1, x_n)} = 0$; $\hat{\theta}(x_1, x_n) = \bar{x}$ - the unique solution. Verify that $\frac{\partial S(\theta/x_1, x_n)}{\partial \theta}|_{\theta=\bar{x}} = -\frac{n}{\sigma_0^2} < 0$.

Ex. Exponential model

Let $X_1, \dots, X_n \sim \text{Exp}(\frac{1}{\theta})$ for $\theta > 0$.

$$\Rightarrow L(\theta | X_1, \dots, X_n) = \frac{1}{\theta^n} e^{-\frac{n\bar{x}}{\theta}} \Rightarrow \ell(\theta | X_1, \dots, X_n) = -n \ln \theta - \frac{n\bar{x}}{\theta},$$

$$\text{and } S(\theta | X_1, \dots, X_n) = -\frac{1}{\theta} + \frac{n\bar{x}}{\theta^2} / = 0,$$

$$\text{get } \hat{\theta}(X_1, \dots, X_n) = \bar{x} \text{ since } \bar{x} > 0,$$

$$\frac{\partial S(\theta | X_1, \dots, X_n)}{\partial \theta} /_{\theta=\bar{x}} = \frac{1}{\bar{x}^2} - 2 \frac{n\bar{x}}{\bar{x}^3} /_{\theta=\bar{x}} = -\frac{1}{\bar{x}^2} < 0;$$

so \bar{x} is the MLE.

Ex. (No simple formulas for MLE):

Three categories 1, 2, 3. Probabilities: $p_1 = \theta$, $p_2 = \theta^2$,

$p_3 = 1 - \theta - \theta^2$ where $\theta \in [0, (\sqrt{5}-1)/2] = [0, 0.61803]$ is unknown.

Note that $\theta + \theta^2 \leq 1$ imposes the upper bound on θ ; obtain this bound using $\theta + \theta^2 - 1 = 0$ for θ . This relationship between p_i 's arises in genetics.

For a sample of n (small relative to the population size), so if x_i are iid, $L(\theta | X_1, \dots, X_n) = \theta^{x_1} \theta^{2x_2} (1 - \theta - \theta^2)^{x_3}$, x_i - sample count in the i th class.

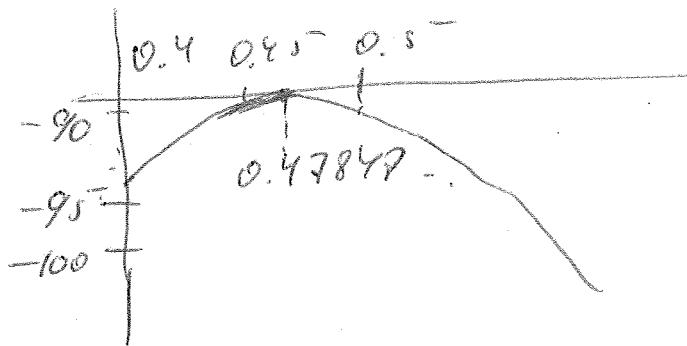
$$\Rightarrow S(\theta | s_1, \dots, s_n) = \frac{x_1 + 2x_2}{\theta} - \frac{x_3(1+2\theta)}{1-\theta-\theta^2};$$

$$\Rightarrow \hat{\theta} \text{ is the root of the quadratic equation}$$

$$(x_1 + 2x_2)(1 - \theta - \theta^2) - x_3(\theta + 2\theta^2) = -(x_1 + 2x_2 + 2x_3)\theta^2 - (x_1 + 2x_2 + x_3)\theta + (x_1 + 2x_2)$$

the resulting $\hat{\theta}$ depends on x_1, x_2, x_3 - it is not even guaranteed to be in $[0, 1]$

in general! 4 possible values - 2 roots, & boundary points 0 and 0.61803. Typically, we have to evaluate the numerically by evaluating the likelihood at two points. So if $x_1 = 70$, $x_2 = 5$, $x_3 = 25$, ~~$\hat{\theta}_1 = -1.28616$~~ , $\hat{\theta}_2 = 0.47849$,

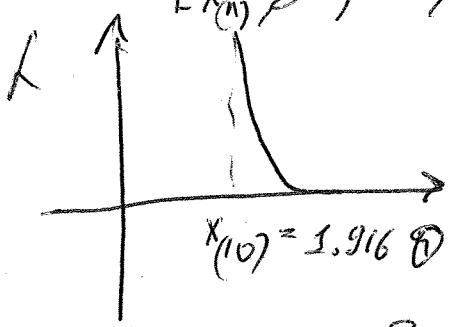


Ex. 6.2.5 Uniform $[0, \theta]$ model - attention to smoothness conditions.

Let $x_1, \dots, x_n \sim \text{Uniform}[0, \theta]$; $\theta > 0$ is unknown.

$$\Rightarrow L(\theta | x_1, \dots, x_n) = \begin{cases} \theta^{-n}, & x_i \leq \theta, i=1, \dots, n \\ 0, & x_i > \theta \text{ for at least one } i \end{cases}$$

$$= \theta^{-n} \Gamma(x_{(n)})^{-1} (\theta); \text{ e.g. plot when } n=10, x_{(n)}=1.916$$



$\Rightarrow \hat{\theta} = x_{(n)}$ - cannot be obtained by differentiation!

Ex. 6.2.6 Location-Scale Normal Model

$(x_1, \dots, x_n) \sim N(\mu, \sigma^2)$; $\mu, \sigma^2 > 0$ are unknown;

$$\theta = (\mu, \sigma^2) \in \mathcal{S} = \mathbb{R}^2 \times (0, \infty) ; \quad -\frac{n}{2} \log \sigma^2 - \frac{n}{2\sigma^2} (\bar{x} - \mu)^2 - \frac{n-1}{2\sigma^2} s^2 ;$$

$$\Rightarrow L(\mu, \sigma^2 | x_1, \dots, x_n) = (2\pi\sigma^2)^{-n/2} e^{-\frac{n}{2\sigma^2} (\bar{x} - \mu)^2} e^{-\frac{n-1}{2\sigma^2} s^2}$$

$$\Rightarrow \ell(\mu, \sigma^2 | x_1, \dots, x_n) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \frac{n}{2\sigma^2} (\bar{x} - \mu)^2 - \frac{n-1}{2\sigma^2} s^2$$

First, $\mu = \bar{x}$; plugging it in, we get $-\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \frac{n-1}{2\sigma^2} s^2$

$$\Rightarrow \frac{\partial}{\partial \sigma^2} : -\frac{n}{2\sigma^2} + \frac{n-1}{2(\sigma^2)^2} s^2 = 0$$

The final MLE of (μ, σ^2) is $\left(\bar{x}, \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2\right)$.

Often, we need to estimate some $\psi(\theta)$. Can we say that $\hat{\psi}(\hat{\theta})$ is the MLE of $\psi(\theta)$? If ψ is 1-1, then yes!

If ψ is not 1-1, we can often find a complementary function λ on \mathcal{Q} s.t. (ψ, λ) is a 1-1 function of ψ .

Then, $(\hat{\psi}(s), \hat{\lambda}(s)) = (\psi(\hat{\theta}(s)), \lambda(\hat{\theta}(s)))$ is the joint MLE but ~~$\hat{\psi}(s)$~~ $\hat{\psi}(s)$ is technically not an MLE. Taken separately, it may perform quite poorly.

Ex. 6.2.7. Let $X_i \sim N(\mu_i, 1)$, $i=1, \dots, n$ - independent & with unknown means. Thus, $\theta = (\mu_1, \dots, \mu_n)$, $\Omega = \mathbb{R}^n$. Need to estimate $\psi(\theta) = \mu_1^2 + \dots + \mu_n^2$.

$$\Rightarrow L(\theta | x_1, \dots, x_n) = -\frac{1}{2} \sum_{i=1}^n (x_i - \mu_i)^2, \text{ have } \hat{\theta}(x_1, \dots, x_n) = (x_1, \dots, x_n). \text{ The "plug-in" MLE is } \hat{\psi} = \sum_{i=1}^n x_i^2$$

$$\text{However, } E_\theta \left(\sum_{i=1}^n x_i^2 \right) = \sum_{i=1}^n [E_\theta(x_i^2) + V(x_i)] = n + \psi(\theta)$$

; here $E_\theta g$ refers to the expectation of $g(s)$ when $s \sim \theta$.

When n is large, $\hat{\psi}$ is clearly far from the true value
; should be using $\sum_{i=1}^n x_i^2 - n$ instead.