

# Maximum Likelihood Estimation

Assume  $\mathbf{X} \sim P_\theta, \theta \in \Theta$ , with joint pdf (or pmf)  $f(\mathbf{x} | \theta)$ .

Suppose we observe  $\mathbf{X} = \mathbf{x}$ .

The **Likelihood function** is

$$L(\theta | \mathbf{x}) = f(\mathbf{x} | \theta)$$

as a function of  $\theta$  (with the data  $\mathbf{x}$  held fixed).

The likelihood function  $L(\theta | \mathbf{x})$  and joint pdf  $f(\mathbf{x} | \theta)$  are the same except that  $f(\mathbf{x} | \theta)$  is generally viewed as a function of  $\mathbf{x}$  with  $\theta$  held fixed, and  $L(\theta | \mathbf{x})$  as a function of  $\theta$  with  $\mathbf{x}$  held fixed.

$f(\mathbf{x} | \theta)$  is a density in  $\mathbf{x}$  for each fixed  $\theta$ .

But  $L(\theta | \mathbf{x})$  is **not** a density (or mass function) in  $\theta$  for fixed  $\mathbf{x}$  (except by coincidence).

## The Maximum Likelihood Estimator (MLE)

A point estimator  $\hat{\theta} = \hat{\theta}(\mathbf{x})$  is a MLE for  $\theta$  if

$$L(\hat{\theta} | \mathbf{x}) = \sup_{\theta} L(\theta | \mathbf{x}),$$

that is,  $\hat{\theta}$  maximizes the likelihood.

In most cases, the maximum is achieved at a unique value, and we can refer to “the” MLE, and write

$$\hat{\theta}(\mathbf{x}) = \operatorname{argmax}_{\theta} L(\theta | \mathbf{x}).$$

(But there are cases where the likelihood has flat spots and the MLE is not unique.)

## Motivation for MLE's

Note: We often write  $L(\theta | \mathbf{x}) = L(\theta)$ , suppressing  $\mathbf{x}$ , which is kept fixed at the observed data.

Suppose  $\mathbf{x} \in \mathbb{R}^n$ .

Discrete Case:

If  $f(\cdot | \theta)$  is a mass function ( $\mathbf{X}$  is discrete), then

$$L(\theta) = f(\mathbf{x} | \theta) = P_\theta(\mathbf{X} = \mathbf{x}).$$

$L(\theta)$  is the probability of getting the observed data  $\mathbf{x}$  when the parameter value is  $\theta$ .

Continuous Case:

When  $f(\cdot | \theta)$  is a continuous density  $P_\theta(\mathbf{X} = \mathbf{x}) = 0$ , but if  $B \subset \mathbb{R}^n$  is a very, very small ball (or cube) centered at the observed data  $\mathbf{x}$ , then

$$P_\theta(\mathbf{X} \in B) \approx f(\mathbf{x} | \theta) \times \text{Volume}(B) \propto L(\theta).$$

$L(\theta)$  is proportional to the probability the random data  $\mathbf{X}$  will be **close** to the observed data  $\mathbf{x}$  when the parameter value is  $\theta$ .

Thus, the MLE  $\hat{\theta}$  is the value of  $\theta$  which makes the observed data  $\mathbf{x}$  “most probable”.

To find  $\hat{\theta}$ , we maximize  $L(\theta)$ . This is usually done by calculus (finding a stationary point), but **not** always.

If the parameter space  $\Theta$  contains endpoints or boundary points, the maximum can be achieved at a boundary point without being a stationary point.

If  $L(\theta)$  is not “smooth” (continuous and everywhere differentiable), the maximum does **not** have to be achieved at a stationary point.

### Cautionary Example:

Suppose  $X_1, \dots, X_n$  are iid  $\text{Uniform}(0, \theta)$  and  $\Theta = (0, \infty)$ .

Given data  $\mathbf{x} = (x_1, \dots, x_n)$ , find the MLE for  $\theta$ .

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n \theta^{-1} I(0 \leq x_i \leq \theta) = \theta^{-n} I(0 \leq \min x_i) I(\max x_i \leq \theta) \\ &= \begin{cases} \theta^{-n} & \text{for } \theta \geq \max x_i \\ 0 & \text{for } 0 < \theta < \max x_i \end{cases} \quad (\text{Draw this!}) \end{aligned}$$

which is maximized at  $\theta = \max x_i$ , which is a point of discontinuity (and certainly **not** a stationary point).

Thus, the MLE is  $\hat{\theta} = \max x_i = x_{(n)}$ .

Notes:

$L(\theta) = 0$  for  $\theta < \max x_i$  is just saying that these values of  $\theta$  are absolutely ruled out by the data (which is obvious).

A strange property of the MLE in this example (not typical):

$$P_{\theta}(\hat{\theta} < \theta) = 1$$

The MLE is biased; it is always less than the true value.

### A Similar Example:

Let  $X_1, \dots, X_n$  be iid  $\text{Uniform}(\alpha, \beta)$  and  $\Theta = \{(\alpha, \beta) : \alpha < \beta\}$ .

Given data  $\mathbf{x} = (x_1, \dots, x_n)$ , find the MLE for  $\theta = (\alpha, \beta)$ .

$$\begin{aligned} L(\alpha, \beta) &= \prod_{i=1}^n (\beta - \alpha)^{-1} I(\alpha \leq x_i \leq \beta) \\ &= (\beta - \alpha)^{-n} I(\alpha \leq \min x_i) I(\max x_i \leq \beta) \\ &= \begin{cases} (\beta - \alpha)^{-n} & \text{for } \alpha \leq \min x_i, \max x_i \leq \beta \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

which is maximized by making  $\beta - \alpha$  as small as possible without entering “0 otherwise” region.

Clearly, the maximum is achieved at  $(\alpha, \beta) = (\min x_i, \max x_i)$ . Thus the MLE is  $\hat{\theta} = (\hat{\alpha}, \hat{\beta}) = (\min x_i, \max x_i)$ .

Again,  $P_{\alpha, \beta}(\alpha < \hat{\alpha}, \hat{\beta} < \beta) = 1$ .

## Maximizing the Likelihood (one parameter)

**Basic Result:** A continuous function  $g(\theta)$  defined on a closed, bounded interval  $J$  attains its supremum (but might do so at one of the endpoints).

(That is, there exists a point  $\theta_0 \in J$  such that  $g(\theta_0) = \sup_{\theta \in J} g(\theta)$ . )

**Consequence:** Suppose  $g(\theta)$  is a continuous, non-negative function defined on an open interval  $J = (c, d)$  (where perhaps  $c = -\infty$  or  $d = +\infty$ ). If

$$\lim_{\theta \rightarrow c} g(\theta) = \lim_{\theta \rightarrow d} g(\theta) = 0,$$

then  $g$  attains its supremum.

- Thus, MLE's usually exist when the likelihood function is continuous.

Suppose the function  $g(\theta)$  is defined on an interval  $\Theta$  (which may be open or closed, infinite or finite).

If  $g$  is differentiable and attains its supremum at a point  $\theta_0$  in the interior of  $\Theta$ , that point must be a stationary point (that is,  $g'(\theta_0) = 0$ ).

(1) If  $g'(\theta_0) = 0$  and  $g''(\theta_0) < 0$ , then  $\theta_0$  is a local maximum (but might not be the global maximum).

(2) If  $g'(\theta_0) = 0$  and  $g''(\theta) < 0$  for **all**  $\theta \in \Theta$ , then  $\theta_0$  is a global maximum (that is, it attains the supremum).

(1) is necessary (but not sufficient) for  $\theta_0$  to be a global maximum. (2) is sufficient (but not necessary).

A function satisfying  $g''(\theta) < 0$  for all  $\theta \in \Theta$  is called **strictly concave**. It lies below any tangent line.

## Maximizing the Likelihood (multi-parameter)

**Basic Result:** A continuous function  $g(\theta)$  defined on a closed, bounded set  $J \subset R^k$  attains its supremum (but might do so on the boundary).

**Consequence:** Suppose  $g(\theta)$  is a continuous, non-negative function defined for all  $\theta \in R^k$ . If  $g(\theta) \rightarrow 0$  as  $\theta \rightarrow \infty$ , then  $g$  attains its supremum.

- Thus, MLE's usually exist when the likelihood function is continuous.

Suppose the function  $g(\theta)$  is defined on a **convex** set  $\Theta \subset R^k$  (that is, the line segment joining any two points in  $\Theta$  lies entirely inside  $\Theta$ ).

If  $g$  is differentiable and attains its supremum at a point  $\theta_0$  in the interior of  $\Theta$ , that point must be a stationary point:

$$\frac{\partial g(\theta_0)}{\partial \theta_i} = 0 \quad \text{for } i = 1, 2, \dots, k.$$

Define the gradient vector  $D$  and Hessian matrix  $H$ :

$$D(\theta) = \left( \frac{\partial g(\theta)}{\partial \theta_i} \right)_{i=1}^k \quad (\text{a } k \times 1 \text{ vector}).$$

$$H(\theta) = \left( \frac{\partial^2 g(\theta)}{\partial \theta_i \partial \theta_j} \right)_{i,j=1}^k \quad (\text{a } k \times k \text{ matrix}).$$

where  $\theta = (\theta_1, \theta_2, \dots, \theta_k)'$ .

(1) If  $D(\theta_0) = 0$  and  $H(\theta_0)$  is **negative definite**, then  $\theta_0$  is a local maximum (but might not be the global maximum).

(2) If  $D(\theta_0) = 0$  and  $H(\theta)$  is negative definite for **all**  $\theta \in \Theta$ , then  $\theta_0$  is a global maximum (that is, it attains the supremum).

(1) is necessary (but not sufficient) for  $\theta_0$  to be a global maximum. (2) is sufficient (but not necessary).

A function for which  $H(\theta)$  is negative definite for all  $\theta \in \Theta$  is called **strictly concave**. It lies below any tangent plane.

### Example:

Observe  $X_1, \dots, X_n$  be iid  $\text{Gamma}(\alpha, \beta)$ .

### Preliminaries:

$$(\text{likelihood}) \quad L(\alpha, \beta) = \prod_{i=1}^n \frac{x_i^{\alpha-1} e^{-x_i/\beta}}{\beta^\alpha \Gamma(\alpha)}.$$

Maximizing  $L$  is same as maximizing  $\ell = \log L$  given by

$$\begin{aligned} \ell(\alpha, \beta) &= (\alpha - 1)T_1 - T_2/\beta - n\alpha \log \beta - n \log \Gamma(\alpha) \\ \text{where } T_1 &= \sum_i \log x_i, \quad T_2 = \sum_i x_i. \end{aligned}$$

Note that  $T = (T_1, T_2)$  is the natural sufficient statistic of this 2pef.

$$\frac{\partial \ell}{\partial \alpha} = T_1 - n \log \beta - n\psi(\alpha)$$

$$\text{where } \psi(\alpha) \equiv \frac{d}{d\alpha} \log \Gamma(\alpha) = \frac{\Gamma'(\alpha)}{\Gamma(\alpha)}$$

$$\frac{\partial \ell}{\partial \beta} = \frac{T_2}{\beta^2} - \frac{n\alpha}{\beta} = \frac{1}{\beta^2} (T_2 - n\alpha\beta)$$

$$\frac{\partial^2 \ell}{\partial \alpha^2} = -n\psi'(\alpha)$$

$$\frac{\partial^2 \ell}{\partial \beta^2} = \frac{-2T_2}{\beta^3} + \frac{n\alpha}{\beta^2} = \frac{-1}{\beta^3} (2T_2 - n\alpha\beta)$$

$$\frac{\partial^2 \ell}{\partial \alpha \partial \beta} = \frac{-n}{\beta}$$



**Situation #1:** Suppose  $\alpha = \alpha_0$  is known. Find MLE for  $\beta$ .

(Drop  $\alpha$  from arguments:  $\ell(\beta) = \ell(\alpha_0, \beta)$  etc.)

$\ell(\beta)$  is continuous and differentiable.

$\ell(\beta)$  has a unique stationary point:

$$\begin{aligned}\ell'(\beta) &= \frac{\partial \ell}{\partial \beta} = \frac{1}{\beta^2} (T_2 - n\alpha_0\beta) = 0 \\ \text{iff } T_2 &= n\alpha_0\beta \quad \text{iff } \beta = \frac{T_2}{n\alpha_0} (\equiv \beta^*).\end{aligned}$$

Now we check the second derivative.

$$\ell''(\beta) = \frac{\partial^2 \ell}{\partial \beta^2} = \frac{-1}{\beta^3} (2T_2 - n\alpha\beta) = \frac{-1}{\beta^3} (T_2 + (T_2 - n\alpha\beta)).$$

Note  $\ell''(\beta^*) < 0$  since  $T_2 - n\alpha_0\beta^* = 0$ , but  $\ell''(\beta) > 0$  for  $\beta > 2T_2/(n\alpha_0)$ .

Thus, the stationary point satisfies the necessary condition for a global maximum, but **not** the sufficient condition (i.e.,  $\ell(\beta)$  is **not** a strictly concave function).

How can we be sure that we have found the global maximum, and not just a local maximum?

In this case, there is a simple argument: The stationary point  $\beta^*$  is unique, and  $\ell'(\beta) > 0$  for  $\beta < \beta^*$ , and  $\ell'(\beta) < 0$  for  $\beta > \beta^*$ . This ensures  $\beta^*$  is the unique global maximizer.

Conclusion:  $\hat{\beta} = \frac{T_2}{n\alpha_0}$  is the MLE.

(This is a function of  $T_2$ , which is a sufficient statistic for  $\beta$  when  $\alpha$  is known.)

**Situation #2:** Suppose  $\beta = \beta_0$  is known. Find MLE for  $\alpha$ .

(Drop  $\beta$  from arguments:  $\ell(\alpha) = \ell(\alpha, \beta_0)$  etc.)

Note:  $\ell'(\alpha)$  and  $\ell''(\alpha)$  involve  $\psi(\alpha)$  The function  $\psi$  is infinitely differentiable on the interval  $(0, \infty)$ , and satisfies  $\psi'(\alpha) > 0$  and  $\psi''(\alpha) < 0$  for all  $\alpha > 0$ . (The function is strictly increasing and strictly concave.) Also

$$\lim_{\alpha \rightarrow 0^+} \psi(\alpha) = -\infty \quad \text{and} \quad \lim_{\alpha \rightarrow \infty} \psi(\alpha) = \infty. \quad (\text{Draw a Picture.})$$

Thus  $\psi^{-1} : \mathbb{R} \rightarrow (0, \infty)$  exists.

$\ell(\alpha)$  is continuous and differentiable.

$\ell(\alpha)$  has a unique stationary point:

$$\begin{aligned} \ell'(\alpha) &= T_1 - n \log \beta_0 - n\psi(\alpha) = 0 \\ \text{iff } \psi(\alpha) &= T_1/n - \log \beta_0 \\ \text{iff } \alpha &= \psi^{-1}(T_1/n - \log \beta_0) \end{aligned}$$

This is the unique global maximizer since

$$\ell''(\alpha) = -n\psi'(\alpha) < 0 \quad \text{for all } \alpha > 0.$$

Thus  $\hat{\alpha} = \psi^{-1}(T_1/n - \log \beta_0)$  is the MLE.

(This is a function of  $T_1$ , which is a sufficient statistic for  $\alpha$  when  $\beta$  is known.)

**Situation #3:** Find MLE for  $\theta = (\alpha, \beta)$

$\ell(\alpha, \beta)$  is continuous and differentiable.

A stationary point must satisfy the system of two equations:

$$\begin{aligned}\frac{\partial \ell}{\partial \alpha} &= T_1 - n \log \beta - n\psi(\alpha) = 0 \\ \frac{\partial \ell}{\partial \beta} &= \frac{1}{\beta^2} (T_2 - n\alpha\beta) = 0\end{aligned}$$

Solving the second equation for  $\beta$  gives

$$\beta = \frac{T_2}{n\alpha}.$$

Plugging this into the first equation, and rearranging a bit leads to

$$\frac{T_1}{n} - \log \left( \frac{T_2}{n} \right) = \psi(\alpha) - \log \alpha \equiv H(\alpha)$$

The function  $H(\alpha)$  is continuous and strictly increasing from  $(0, \infty)$  to  $(-\infty, 0)$ , so that it has an inverse mapping  $(-\infty, 0)$  to  $(0, \infty)$ .

Thus, the solution to the above equation can be written:

$$\alpha = H^{-1} \left( \frac{T_1}{n} - \log \left( \frac{T_2}{n} \right) \right).$$

Thus, the unique stationary point is:

$$\begin{aligned}\hat{\alpha} &= H^{-1} \left( \frac{T_1}{n} - \log \left( \frac{T_2}{n} \right) \right) \\ \hat{\beta} &= \frac{T_2}{n\hat{\alpha}}.\end{aligned}$$

Is this the MLE?

Let us examine the Hessian.