# 2 Maximum Likelihood Estimation (MLE,       ) — More Formally Review

1. We have random variables $X_1$, $X_2$, $\cdots$, $X_n$, which are assumed to be mutually independently and identically distributed.

2. The distribution function of $\{X_i\}_{i=1}^n$ is $f(x;\theta)$, where $x = (x_1, x_2, \cdots, x_n)$ and $\theta = (\mu, \Sigma)$.

   Note that $X$ is a vector of random variables and $x$ is a vector of their realizations (i.e., observed data).

   Likelihood function $L(\cdot)$ is defined as $L(\theta; x) = f(x; \theta)$.

   Note that $f(x; \theta) = \prod_{i=1}^n f(x_i; \theta)$ when $X_1$, $X_2$, $\cdots$, $X_n$ are mutually indepen-

dently and identically distributed.

The maximum likelihood estimator (MLE) of $\theta$ is $\theta$ such that:

$$\max_{\theta}\ L(\theta; X). \qquad \Longleftrightarrow \qquad \max_{\theta}\ \log L(\theta; X).$$

MLE satisfies the following two conditions:

(a) $\dfrac{\partial \log L(\theta; X)}{\partial \theta} = 0.$

(b) $\dfrac{\partial^2 \log L(\theta; X)}{\partial \theta \partial \theta'}$ is a negative definite matrix.

3. **Fisher's information matrix (                              )** is defined as:

$$I(\theta) = -\mathrm{E}\Big(\frac{\partial^2 \log L(\theta; X)}{\partial \theta \partial \theta'}\Big),$$

where we have the following equality:

$$-\mathrm{E}\Big(\frac{\partial^2 \log L(\theta; X)}{\partial \theta \partial \theta'}\Big) = \mathrm{E}\Big(\frac{\partial \log L(\theta; X)}{\partial \theta}\frac{\partial \log L(\theta; X)}{\partial \theta'}\Big) = \mathrm{V}\Big(\frac{\partial \log L(\theta; X)}{\partial \theta}\Big)$$

**Proof of the above equality:**

$$\int L(\theta; x)\mathrm{d}x = 1$$

Take a derivative with respect to $\theta$.

$$\int \frac{\partial L(\theta; x)}{\partial \theta}\mathrm{d}x = 0$$

(We assume that (i) the domain of $x$ does not depend on $\theta$ and (ii) the derivative $\frac{\partial L(\theta; x)}{\partial \theta}$ exists.)

Rewriting the above equation, we obtain:

$$\int \frac{\partial \log L(\theta; x)}{\partial \theta} L(\theta; x)\mathrm{d}x = 0,$$

i.e.,

$$\mathrm{E}\left(\frac{\partial \log L(\theta; X)}{\partial \theta}\right) = 0.$$

28

Again, differentiating the above with respect to $\theta$, we obtain:

$$
\int \frac{\partial^2 \log L(\theta; x)}{\partial \theta \partial \theta'} L(\theta; x) \mathrm{d}x + \int \frac{\partial \log L(\theta; x)}{\partial \theta} \frac{\partial L(\theta; x)}{\partial' \theta} \mathrm{d}x
$$
$$
= \int \frac{\partial^2 \log L(\theta; x)}{\partial \theta \partial \theta'} L(\theta; x) \mathrm{d}x + \int \frac{\partial \log L(\theta; x)}{\partial \theta} \frac{\partial \log L(\theta; x)}{\partial \theta'} L(\theta; x) \mathrm{d}x
$$
$$
= \mathrm{E}\Big(\frac{\partial^2 \log L(\theta; X)}{\partial \theta \partial \theta'}\Big) + \mathrm{E}\Big(\frac{\partial \log L(\theta; X)}{\partial \theta} \frac{\partial \log L(\theta; X)}{\partial \theta'}\Big) = 0.
$$

Therefore, we can derive the following equality:

$$
-\mathrm{E}\left(\frac{\partial^2 \log L(\theta; X)}{\partial \theta \partial \theta'}\right) = \mathrm{E}\left(\frac{\partial \log L(\theta; X)}{\partial \theta} \frac{\partial \log L(\theta; X)}{\partial \theta'}\right) = \mathrm{V}\left(\frac{\partial \log L(\theta; X)}{\partial \theta}\right),
$$

where the second equality utilizes $\mathrm{E}\left(\dfrac{\partial \log L(\theta; X)}{\partial \theta}\right) = 0$.

4. **Cramer-Rao Lower Bound (                                    ):** $(I(\theta))^{-1}$

Suppose that an unbiased estimator of $\theta$ is given by $s(X)$.

Then, we have the following:

$$V(s(X)) \geq (I(\theta))^{-1}$$

**Proof:**

The expectation of $s(X)$ is:

$$E(s(X)) = \int s(x)L(\theta; x)dx.$$

Differentiating the above with respect to $\theta$,

$$\frac{\partial E(s(X))}{\partial \theta'} = \int s(x)\frac{\partial L(\theta; x)}{\partial \theta'}dx = \int s(x)\frac{\partial \log L(\theta; x)}{\partial \theta'}L(\theta; x)dx$$
$$= \text{Cov}\left(s(X), \frac{\partial \log L(\theta; X)}{\partial \theta}\right)$$

30

For simplicity, let $s(X)$ and $\theta$ be scalars.

Then,

$$\left(\frac{\partial \mathrm{E}(s(X))}{\partial \theta}\right)^2 = \left(\mathrm{Cov}\left(s(X), \frac{\partial \log L(\theta; X)}{\partial \theta}\right)\right)^2 = \rho^2 \mathrm{V}\left(s(X)\right) \mathrm{V}\left(\frac{\partial \log L(\theta; X)}{\partial \theta}\right)$$
$$\leq \mathrm{V}\left(s(X)\right) \mathrm{V}\left(\frac{\partial \log L(\theta; X)}{\partial \theta}\right),$$

where $\rho$ denotes the correlation coefficient between $s(X)$ and $\dfrac{\partial \log L(\theta; X)}{\partial \theta}$, i.e.,

$$\rho = \frac{\mathrm{Cov}\left(s(X), \dfrac{\partial \log L(\theta; X)}{\partial \theta}\right)}{\sqrt{\mathrm{V}\left(s(X)\right)}\sqrt{\mathrm{V}\left(\dfrac{\partial \log L(\theta; X)}{\partial \theta}\right)}}.$$

Note that $|\rho| \leq 1$.

Therefore, we have the following inequality:

$$\left(\frac{\partial \mathrm{E}(s(X))}{\partial \theta}\right)^2 \leq \mathrm{V}(s(X)) \, \mathrm{V}\left(\frac{\partial \log L(\theta; X)}{\partial \theta}\right),$$

i.e.,

$$\mathrm{V}(s(X)) \geq \frac{\left(\dfrac{\partial \mathrm{E}(s(X))}{\partial \theta}\right)^2}{\mathrm{V}\left(\dfrac{\partial \log L(\theta; X)}{\partial \theta}\right)}$$

Especially, when $\mathrm{E}(s(X)) = \theta$,

$$\mathrm{V}(s(X)) \geq \frac{1}{-\mathrm{E}\left(\dfrac{\partial^2 \log L(\theta; X)}{\partial \theta^2}\right)} = (I(\theta))^{-1}.$$

Even in the case where $s(X)$ is a vector, the following inequality holds.

$$\mathrm{V}(s(X)) \geq (I(\theta))^{-1},$$

32

where $I(\theta)$ is defined as:

$$
\begin{aligned}
I(\theta) &= -\mathrm{E}\left(\frac{\partial^2 \log L(\theta; X)}{\partial\theta\partial\theta'}\right) \\
&= \mathrm{E}\left(\frac{\partial \log L(\theta; X)}{\partial\theta}\frac{\partial \log L(\theta; X)}{\partial\theta'}\right) = \mathrm{V}\left(\frac{\partial \log L(\theta; X)}{\partial\theta}\right).
\end{aligned}
$$

The variance of any unbiased estimator of $\theta$ is larger than or equal to $(I(\theta))^{-1}$.

5. Asymptotic Normality of MLE:

Let $\tilde{\theta}$ be MLE of $\theta$.

As $n$ goes to infinity, we have the following result:

$$\sqrt{n}(\tilde{\theta} - \theta) \longrightarrow N\left(0, \lim_{n\to\infty}\left(\frac{I(\theta)}{n}\right)^{-1}\right),$$

where it is assumed that $\lim_{n\to\infty}\left(\dfrac{I(\theta)}{n}\right)$ converges.

That is, when $n$ is large, $\tilde{\theta}$ is approximately distributed as follows:

$$\tilde{\theta} \sim N\left(\theta, (I(\theta))^{-1}\right).$$

Suppose that $s(X) = \tilde{\theta}$.

When $n$ is large, $V(s(X))$ is approximately equal to $(I(\theta))^{-1}$.

Practically, we utilize the following approximated distribution:

$$\tilde{\theta} \sim N\left(\theta, (I(\tilde{\theta}))^{-1}\right).$$

Then, we can obtain the significance test and the confidence interval for $\theta$

6. **Central Limit Theorem:** Let $X_1$, $X_2$, $\cdots$, $X_n$ be mutually independently distributed random variables with mean $E(X_i) = \mu$ and variance $V(X_i) = \sigma^2 < \infty$ for $i = 1, 2, \cdots, n$.

Define $\overline{X} = (1/n) \sum_{i=1}^{n} X_i$.

Then, the central limit theorem is given by:

$$\frac{\overline{X} - E(\overline{X})}{\sqrt{V(\overline{X})}} = \frac{\overline{X} - \mu}{\sigma / \sqrt{n}} \longrightarrow N(0, 1).$$

Note that $E(\overline{X}) = \mu$ and $V(\overline{X}) = \sigma^2/n$.

That is,

$$\sqrt{n}(\overline{X} - \mu) = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} (X_i - \mu) \longrightarrow N(0, \sigma^2).$$

Note that $E(\overline{X}) = \mu$ and $nV(\overline{X}) = \sigma^2$.

In the case where $X_i$ is a vector of random variable with mean $\mu$ and variance $\Sigma < \infty$, the central limit theorem is given by:

$$\sqrt{n}(\overline{X} - \mu) = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} (X_i - \mu) \longrightarrow N(0, \Sigma).$$

Note that $E(\overline{X}) = \mu$ and $nV(\overline{X}) = \Sigma$.

7. **Central Limit Theorem II:** Let $X_1$, $X_2$, $\cdots$, $X_n$ be mutually independently distributed random variables with mean $E(X_i) = \mu$ and variance $V(X_i) = \sigma_i^2$ for $i = 1, 2, \cdots, n$.

Assume:

$$\sigma^2 = \lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} \sigma_i^2 < \infty.$$

Define $\overline{X} = (1/n) \sum_{i=1}^{n} X_i$.

Then, the central limit theorem is given by:

$$\frac{\overline{X} - E(\overline{X})}{\sqrt{V(\overline{X})}} = \frac{\overline{X} - \mu}{\sigma/\sqrt{n}} \longrightarrow N(0, 1),$$

i.e.,

$$\sqrt{n}(\overline{X} - \mu) = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} (X_i - \mu) \longrightarrow N(0, \sigma^2).$$

Note that $E(\overline{X}) = \mu$ and $nV(\overline{X}) \longrightarrow \sigma^2$.

37

In the case where $X_i$ is a vector of random variable with mean $\mu$ and variance $\Sigma_i$, the central limit theorem is given by:

$$\sqrt{n}(\overline{X} - \mu) = \frac{1}{\sqrt{n}} \sum_{i=1}^{n}(X_i - \mu) \longrightarrow N(0, \Sigma),$$

where $\Sigma = \lim_{n \to \infty} \dfrac{1}{n} \sum_{i=1}^{n} \Sigma_i < \infty$.

Note that $E(\overline{X}) = \mu$ and $nV(\overline{X}) \longrightarrow \Sigma$.

[Review of Asymptotic Theories]

• **Convergence in Probability** (          ) $X_n \longrightarrow a$, i.e., $X$ converges in probability to $a$, where $a$ is a fixed number.

• **Convergence in Distribution** (　　　) $X_n \longrightarrow X$, i.e., $X$ converges in distribution to $X$. The distribution of $X_n$ converges to the distribution of $X$ as $n$ goes to infinity.

**Some Formulas**

$X_n$ and $Y_n$ : Convergence in Probability

$Z_n$ : Convergence in Distribution

• If $X_n \longrightarrow a$, then $f(X_n) \longrightarrow f(a)$.

• If $X_n \longrightarrow a$ and $Y_n \longrightarrow b$, then $f(X_n Y_n) \longrightarrow f(ab)$.

• If $X_n \longrightarrow a$ and $Z_n \longrightarrow Z$, then $X_n Z_n \longrightarrow aZ$, i.e., $aZ$ is distributed with mean $E(aZ) = aE(Z)$ and variance $V(aZ) = a^2 V(Z)$.

**[End of Review]**

8. **Weak Law of Large Numbers (                    ) — Review:**

   $n$ random variables $X_1$, $X_2$, $\cdots$, $X_n$ are assumed to be mutually independently
   and identically distributed, where $E(X_i) = \mu$ and $V(X_i) = \sigma^2 < \infty$.

   Then, $\overline{X} \longrightarrow \mu$ as $n \longrightarrow \infty$, which is called the **weak law of large numbers**.

   $\longrightarrow$ Convergence in probability

   $\longrightarrow$ Proved by Chebyshev's inequality

9. **Some Formulas of Expectaion and Variance in Multivariate Cases**
   **— Review:**

   A vector of randam variavle $X$: $E(X) = \mu$ and $V(X) \equiv E((X - \mu)(X - \mu)') = \Sigma$

   Then, $E(AX) = A\mu$ and $V(AX) = A\Sigma A'$.

**Proof:**

$E(AX) = AE(X) = A\mu$

$V(AX) = E((AX - A\mu)(AX - A\mu)') = E(A(X - \mu)(A(X - \mu))')$

$\quad = E(A(X - \mu)(X - \mu)'A') = AE((X - \mu)(X - \mu)')A' = AV(X)A' = A\Sigma A'$

10. **Asymptotic Normality of MLE — Proof:**

The density (or probability) function of $X_i$ is given by $f(x_i; \theta)$.

The likelihood function is: $L(\theta; x) \equiv f(x; \theta) = \prod_{i=1}^{n} f(x_i; \theta)$,

where $x = (x_1, x_2, \cdots, x_n)$.

MLE of $\theta$ results in the following maximization problem:

$$\max_{\theta} \quad \log L(\theta; x).$$

41

A solution of the above problem is given by MLE of $\theta$, denoted by $\tilde{\theta}$.

That is, $\tilde{\theta}$ is given by the $\theta$ which satisfies the following equation:

$$\frac{\partial \log L(\theta; x)}{\partial \theta} = \sum_{i=1}^{n} \frac{\partial \log f(x_i; \theta)}{\partial \theta} = 0.$$

Replacing $x_i$ by the underlying random variable $X_i$, $\dfrac{\partial \log f(X_i; \theta)}{\partial \theta}$ is taken as the $i$th random variable, i.e., $X_i$ in the **Central Limit Theorem II**.

Consider applying **Central Limit Theorem II** as follows:

$$\frac{\dfrac{1}{n} \sum_{i=1}^{n} \dfrac{\partial \log f(X_i; \theta)}{\partial \theta} - \mathrm{E}\Big(\dfrac{1}{n} \sum_{i=1}^{n} \dfrac{\partial \log f(X_i; \theta)}{\partial \theta}\Big)}{\sqrt{\mathrm{V}\Big(\dfrac{1}{n} \sum_{i=1}^{n} \dfrac{\partial \log f(X_i; \theta)}{\partial \theta}\Big)}} = \frac{\dfrac{1}{n} \dfrac{\partial \log L(\theta; X)}{\partial \theta} - \mathrm{E}\Big(\dfrac{1}{n} \dfrac{\partial \log L(\theta; X)}{\partial \theta}\Big)}{\sqrt{\mathrm{V}\Big(\dfrac{1}{n} \dfrac{\partial \log L(\theta; X)}{\partial \theta}\Big)}}.$$

Note that

$$\sum_{i=1}^{n} \frac{\partial \log f(X_i; \theta)}{\partial \theta} = \frac{\partial \log L(\theta; X)}{\partial \theta}$$

In this case, we need the following expectation and variance:

$$\mathrm{E}\Big(\frac{1}{n}\sum_{i=1}^{n}\frac{\partial \log f(X_i;\theta)}{\partial\theta}\Big) = \mathrm{E}\Big(\frac{1}{n}\frac{\partial \log L(\theta;X)}{\partial\theta}\Big) = 0,$$

and

$$\mathrm{V}\Big(\frac{1}{n}\sum_{i=1}^{n}\frac{\partial \log f(X_i;\theta)}{\partial\theta}\Big) = \mathrm{V}\Big(\frac{1}{n}\frac{\partial \log L(\theta;X)}{\partial\theta}\Big) = \frac{1}{n^2}I(\theta).$$

Note that $\mathrm{E}\Big(\dfrac{\partial \log L(\theta;X)}{\partial\theta}\Big) = 0$ and $\mathrm{V}\Big(\dfrac{\partial \log L(\theta;X)}{\partial\theta}\Big) = I(\theta).$

Thus, the asymptotic distribution of

$$\frac{1}{n}\frac{\partial \log L(\theta; X)}{\partial \theta} = \frac{1}{n}\sum_{i=1}^{n}\frac{\partial \log f(X_i; \theta)}{\partial \theta}$$

is given by:

$$\sqrt{n}\left(\frac{1}{n}\sum_{i=1}^{n}\frac{\partial \log f(X_i; \theta)}{\partial \theta} - \mathrm{E}\Big(\frac{1}{n}\sum_{i=1}^{n}\frac{\partial \log f(X_i; \theta)}{\partial \theta}\Big)\right)$$

$$= \sqrt{n}\left(\frac{1}{n}\frac{\partial \log L(\theta; X)}{\partial \theta} - \mathrm{E}\Big(\frac{1}{n}\frac{\partial \log L(\theta; X)}{\partial \theta}\Big)\right)$$

$$= \frac{1}{\sqrt{n}}\frac{\partial \log L(\theta; X)}{\partial \theta} \longrightarrow N(0, \Sigma)$$

where

$$n\mathrm{V}\Big(\frac{1}{n}\sum_{i=1}^{n}\frac{\partial \log f(X_i; \theta)}{\partial \theta}\Big) = \frac{1}{n}\mathrm{V}\Big(\sum_{i=1}^{n}\frac{\partial \log f(X_i; \theta)}{\partial \theta}\Big) = \frac{1}{n}\mathrm{V}\Big(\frac{\partial \log L(\theta; X)}{\partial \theta}\Big)$$

$$= \frac{1}{n}I(\theta) \longrightarrow \Sigma.$$

44

That is,

$$\frac{1}{\sqrt{n}} \frac{\partial \log L(\theta; X)}{\partial \theta} \ \longrightarrow \ N(0, \Sigma),$$

where $X = (X_1, X_2, \cdots, X_n)$.

Now, replacing $\theta$ by $\tilde{\theta}$, consider the asymptotic distribution of

$$\frac{1}{\sqrt{n}} \frac{\partial \log L(\tilde{\theta}; X)}{\partial \theta},$$

which is expanded around $\tilde{\theta} = \theta$ as follows:

$$0 = \frac{1}{\sqrt{n}} \frac{\partial \log L(\tilde{\theta}; X)}{\partial \theta} \approx \frac{1}{\sqrt{n}} \frac{\partial \log L(\theta; X)}{\partial \theta} + \frac{1}{\sqrt{n}} \frac{\partial^2 \log L(\theta; X)}{\partial \theta \partial \theta'}(\tilde{\theta} - \theta).$$

Therefore,

$$-\frac{1}{\sqrt{n}} \frac{\partial^2 \log L(\theta; X)}{\partial \theta \partial \theta'}(\tilde{\theta} - \theta) \approx \frac{1}{\sqrt{n}} \frac{\partial \log L(\theta; X)}{\partial \theta} \ \longrightarrow \ N(0, \Sigma).$$

45

The left-hand side is rewritten as:

$$-\frac{1}{\sqrt{n}}\frac{\partial^2 \log L(\theta; X)}{\partial\theta\partial\theta'}(\tilde{\theta} - \theta) = \sqrt{n}\left(-\frac{1}{n}\frac{\partial^2 \log L(\theta; X)}{\partial\theta\partial\theta'}\right)(\tilde{\theta} - \theta).$$

Then,

$$\sqrt{n}(\tilde{\theta} - \theta) \approx \left(-\frac{1}{n}\frac{\partial^2 \log L(\theta; X)}{\partial\theta\partial\theta'}\right)^{-1}\left(\frac{1}{\sqrt{n}}\frac{\partial \log L(\theta; X)}{\partial\theta}\right)$$

$$\longrightarrow N(0, \Sigma^{-1}\Sigma\Sigma^{-1}) = N(0, \Sigma^{-1}).$$

Using the law of large number, note that

$$-\frac{1}{n}\frac{\partial^2 \log L(\theta; X)}{\partial\theta\partial\theta'} \longrightarrow \lim_{n\to\infty}\frac{1}{n}\left(-\mathrm{E}\left(\frac{\partial^2 \log L(\theta; X)}{\partial\theta\partial\theta'}\right)\right)$$

$$= \lim_{n\to\infty}\frac{1}{n}\left(\mathrm{V}\left(\frac{\partial \log L(\theta; X)}{\partial}\right)\right) = \lim_{n\to\infty}\frac{1}{n}I(\theta) = \Sigma,$$

46

and $\left(\frac{1}{n}\frac{\partial^2 \log L(\theta; X)}{\partial \theta \partial \theta'}\right)^{-1}\left(\frac{1}{\sqrt{n}}\frac{\partial \log L(\theta; X)}{\partial \theta}\right)$ has the same asymptotic distribution as $\Sigma^{-1}\left(\frac{1}{\sqrt{n}}\frac{\partial \log L(\theta; X)}{\partial \theta}\right)$.

11. **Optimization (       ):**

    MLE of $\theta$ results in the following maximization problem:

    $$\max_{\theta} \ \log L(\theta; x).$$

    We often have the case where the solution of $\theta$ is not derived in closed form.

    $\Longrightarrow$ Optimization procedure

    $$0 = \frac{\partial \log L(\theta; x)}{\partial \theta} = \frac{\partial \log L(\theta^*; x)}{\partial \theta} + \frac{\partial^2 \log L(\theta^*; x)}{\partial \theta \partial \theta'}(\theta - \theta^*).$$

    Solving the above equation with respect to $\theta$, we obtain the following:

    $$\theta = \theta^* - \left(\frac{\partial^2 \log L(\theta^*; x)}{\partial \theta \partial \theta'}\right)^{-1} \frac{\partial \log L(\theta^*; x)}{\partial \theta}.$$

47

Replace the variables as follows:

$$\theta \longrightarrow \theta^{(i+1)}, \qquad \theta^* \longrightarrow \theta^{(i)}.$$

Then, we have:

$$\theta^{(i+1)} = \theta^{(i)} - \left( \frac{\partial^2 \log L(\theta^{(i)}; x)}{\partial \theta \partial \theta'} \right)^{-1} \frac{\partial \log L(\theta^{(i)}; x)}{\partial \theta}.$$

$\Longrightarrow$ **Newton-Raphson method (                    )**

Replacing $\dfrac{\partial^2 \log L(\theta^{(i)}; x)}{\partial \theta \partial \theta'}$ by $\mathrm{E}\left( \dfrac{\partial^2 \log L(\theta^{(i)}; x)}{\partial \theta \partial \theta'} \right)$, we obtain the following optimization algorithm:

$$\theta^{(i+1)} = \theta^{(i)} - \left( \mathrm{E}\left( \frac{\partial^2 \log L(\theta^{(i)}; x)}{\partial \theta \partial \theta'} \right) \right)^{-1} \frac{\partial \log L(\theta^{(i)}; x)}{\partial \theta}$$

$$= \theta^{(i)} + \left( I(\theta^{(i)}) \right)^{-1} \frac{\partial \log L(\theta^{(i)}; x)}{\partial \theta}$$

$\Longrightarrow$ **Method of Scoring (          )**

48