

7 Large Sample Theory

7.1 Convergence

Definition 1 (Convergence in distribution). A sequence of random vectors X_n converges in distribution to X if $P(X_n \leq x) \rightarrow P(X \leq x)$ as $n \rightarrow \infty$ for all points x at which $x \mapsto P(X \leq x)$ is continuous. It is denoted by $X_n \xrightarrow{d} X$ or $X_n \xrightarrow{L} X$.

Example 1. Example: Let X_n be uniform random variable on the set $\{1/n, 2/n, \dots, n/n\}$. Then, for any fixed $x \in [0, 1]$,

$$P(X_n \leq x) = \sum_{i=1}^n \frac{1}{n} \times 1\left(x \leq \frac{i}{n}\right) \rightarrow x.$$

Convergence holds because

$$x = \int_0^x 1 dx = \lim_{n \rightarrow \infty} \sum_{i=1}^n \frac{1}{n} \times 1\left(x \leq \frac{i}{n}\right).$$

Hence, $X_n \xrightarrow{d} X = U(0, 1)$ uniform distribution on $[0, 1]$. So $P(X_n \in A)$ converges to $P(X \in A)$ for all A of the form $\{x : x \leq a\}$.

Remark 1. If the densities converge as $f_n \rightarrow f$, then $P(X_n \in A)$ converges to $P(X \in A)$ for all Borel sets A . But $P(X_n \in A)$ does not converge to $P(X \in A)$ for all A . If $A = \{x : x \text{ is rational}\}$, then $P(X_n \in A) = 1$ but $P(X \in A) = 0$!

Remark 2. Convergence in distribution does not mean that $E[X_n] \rightarrow E[X]$. Suppose that

$$X_n \sim \left(1 - \frac{1}{n}\right) N(0, 1) + \frac{1}{n} N(n^2, 1).$$

Then,

$$P(X_n \leq x) = \left(1 - \frac{1}{n}\right) P\{N(0, 1) \leq x\} + \frac{1}{n} P\{N(n^2, 1) \leq x\} \rightarrow P\{N(0, 1) \leq x\}.$$

But $E[X_n] = n$ and $E[X] = 0$.

Lemma 1 (Portmanteau Lemma). $X_n \xrightarrow{d} X$ if and only if $E[f(X_n)] \rightarrow E[f(X)]$ for all bounded and continuous functions f .

Let $d(x, y)$ be a distance function on \mathbb{R}^k . For example, we can consider the Euclidean distance

$$d(x, y) = \|x - y\|_2 = \sqrt{(x - y)^T (x - y)}.$$

Definition 2 (Convergence in Probability). A sequence of random vectors X_n converges in probability to X if for every $\epsilon > 0$, $P(d(X_n, X) > \epsilon) \rightarrow 0$ (or equivalently $P(d(X_n, X) \leq \epsilon) \rightarrow 1$) as $n \rightarrow \infty$. It is denoted by $X_n \xrightarrow{P} X$.

Example 2. Suppose that X_1, \dots, X_n form an iid sample from some distribution with finite variance. Then,

$$P(|\bar{X} - E[X]| \geq \epsilon) \leq \frac{\text{Var}[\bar{X}]}{\epsilon^2} = \frac{\text{Var}[X]}{n\epsilon^2} \rightarrow 0.$$

Definition 3 (Convergence Almost Surely). A sequence of random vectors X_n converges almost surely to X if

$$P\left(\lim_{n \rightarrow \infty} X_n = X\right) = 1.$$

It is denoted by $X_n \xrightarrow{a.s.} X$. An equivalent definition of convergence almost surely is that $X_n \xrightarrow{a.s.} X$ if and only if, for every $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} P(d(X_k, X) < \epsilon, \text{ for all } k \geq n) = 1.$$

Example 3. Let $Z \sim \text{Uniform}(0, 1)$ and $X_n = 1 (Z < n^{-1})$. Then,

$$\left\{\lim_{n \rightarrow \infty} X_n = 0\right\} = \{Z > 0\}.$$

Hence,

$$P\left(\lim_{n \rightarrow \infty} X_n = 0\right) = P(Z > 0) = 1,$$

that is $X_n \xrightarrow{a.s.} 0$.

Relations among the above types of convergence are included in the following theorem.

Theorem 1. *Basic relationships are as follows.*

1. $X_n \xrightarrow{a.s.} X$ implies $X_n \xrightarrow{P} X$.
2. $X_n \xrightarrow{P} X$ implies $X_n \xrightarrow{d} X$.
3. If $c \in \mathbb{R}^d$ is a constant vector, then $X_n \xrightarrow{d} c$ implies $X_n \xrightarrow{P} c$.

Example 4. Let $X \sim N(0, 1)$. Let $X_1 = X$, $X_2 = -X$, $X_3 = X$, $X_4 = -X$, etc. Then X_i has the same distribution as X for any i . Hence, $X_n \xrightarrow{d} X$. But

$$P(|X_n - X| > 1) = \begin{cases} P(|X - X| > 1) = 0 & n \text{ is odd} \\ P(|-X - X| > 1) = P(|X| > \frac{1}{2}) \approx 0.62. & n \text{ is even} \end{cases}$$

Hence, X_n does not converge in probability to X .

Example 5. Let $Z \sim \text{Uniform}(0, 1)$. Let $X_1 = 1$, $X_2 = 1 (Z < \frac{1}{2})$, $X_3 = 1 (\frac{1}{2} \leq Z < 1)$, $X_4 = 1 (Z < \frac{1}{4})$, $X_5 = 1 (\frac{1}{4} \leq Z < \frac{1}{2})$, In general, if $n = 2^k + m$ for $k \geq 0$ and $0 \leq m < 2^k$, then

$$X_n = 1 \left(\frac{m(n)}{2^{k(n)}} \leq Z < \frac{m(n) + 1}{2^{k(n)}} \right),$$

and $P(X_n = 1) = \frac{1}{2^{k(n)}}$. Hence,

$$P(|X_n - 0| \geq \epsilon) = P(X_n = 1) = \frac{1}{2^{k(n)}} \rightarrow 0.$$

But, for any $Z < 1$, X_n does not converge to any value since we just move the interval $\frac{m}{2^k} \leq Z < \frac{m+1}{2^k}$ from 0 to 1. Hence, we don't have $X_n \xrightarrow{a.s.} 0$.

Theorem 2 (Continuous Mapping Theorem). Let $g : \mathbb{R}^k \mapsto \mathbb{R}^m$ be continuous at every point of a set C such that $P(X \in C) = 1$. Then,

1. If $X_n \xrightarrow{d} X$, then $g(X_n) \xrightarrow{d} g(X)$.
2. If $X_n \xrightarrow{P} X$, then $g(X_n) \xrightarrow{P} g(X)$.
3. If $X_n \xrightarrow{a.s.} X$, then $g(X_n) \xrightarrow{a.s.} g(X)$.

7.2 Consistency of MLE

Let X_1, \dots, X_n be iid with density $p(x | \theta)$. The log-likelihood is

$$\ell(\theta) = \log L(\theta) = \sum_{i=1}^n \log p(x_i | \theta).$$

Lemma 2. Let $\{p(\cdot | \theta) : \theta \in \Theta\}$ be a collection of densities such that the corresponding probability measures satisfy

$$P_\theta \neq P_{\theta_0} \text{ for every } \theta \neq \theta_0 \quad (\text{identification}).$$

Then

$$-KL(p(X | \theta_0), p(X | \theta)) = E \left[\log \left(\frac{p(X | \theta)}{p(X | \theta_0)} \right) \mid \theta_0 \right]$$

attains its maximum uniquely at θ_0 .

Proof. It is a direct consequence of Kullback-Leibler inequality in Estimation chapter. \square

Theorem 3 (Consistency of MLE: Unidimensional). *Suppose that*

1. the distributions P_θ of the observations are distinct,
2. the observations are i.i.d. with probability density $p(\cdot | \theta)$ with respect to some measure μ ,
3. the distributions P_θ have common support so that $\{x : p(x | \theta) > 0\}$ is independent of θ ,
4. the parameter space Θ contains an open set ω of which the true parameter value θ_0 is an interior point,
5. for almost all x , $p(x | \theta)$ is differentiable with respect to θ in ω .

Then, with probability tending to 1 as $n \rightarrow \infty$, the likelihood equation

$$\frac{1}{n} \sum_{i=1}^n \frac{\partial \log p(X_i | \theta)}{\partial \theta} = 0$$

has a root $\hat{\theta}_n$ such that $\hat{\theta}_n$ tends to θ_0 in probability.

Proof. Let $\mathcal{N}(\theta_0, \epsilon)$ be a open interval with the center θ_0 and ϵ as radius. By Assumption 4, we can find a small enough $\epsilon > 0$ such that $\mathcal{N}(\theta_0, \epsilon) \subset \Theta$. Define $x = (x_1, \dots, x_n)$, and

$$S_n = \left\{ x : \sum_{i=1}^n \log p(x_i | \theta_0) > \sum_{i=1}^n \log p(x_i | \theta_0 - \epsilon) \text{ and } \sum_{i=1}^n \log p(x_i | \theta_0) > \sum_{i=1}^n \log p(x_i | \theta_0 + \epsilon) \right\}.$$

By Jensen's inequality,

$$E \left[\log \left(\frac{p(X | \theta)}{p(X | \theta_0)} \right) \mid \theta_0 \right] \leq \log \left(E \left[\frac{p(X | \theta)}{p(X | \theta_0)} \mid \theta_0 \right] \right) = \log \left(\int p(x | \theta) dx \right) = 0,$$

where the equality holds if and only if $\theta = \theta_0$ [Assumption 1]. By LLN,

$$\frac{1}{n} \sum_{i=1}^n \log \left(\frac{p(x_i | \theta)}{p(x_i | \theta_0)} \right) \xrightarrow{P} E \left[\log \left(\frac{p(X | \theta)}{p(X | \theta_0)} \right) \mid \theta_0 \right] < 0, \quad \forall \theta \neq \theta_0,$$

that is, for any $\delta > 0$,

$$P \left(\left| \frac{1}{n} \sum_{i=1}^n \log \left(\frac{p(x_i | \theta)}{p(x_i | \theta_0)} \right) - E \left[\log \left(\frac{p(X | \theta)}{p(X | \theta_0)} \right) \mid \theta_0 \right] \right| < \delta \right) \rightarrow 1.$$

Since the limit is negative, then

$$P\left(\frac{1}{n} \sum_{i=1}^n \log\left(\frac{p(x_i | \theta)}{p(x_i | \theta_0)}\right) < 0\right) = P\left(\sum_{i=1}^n \log p(x_i | \theta_0) > \sum_{i=1}^n \log p(x_i | \theta)\right) \rightarrow 1.$$

Consequently, $P(S_n) \rightarrow 1$. Because $p(x | \theta)$ is differentiable with respect to θ in $\mathcal{N}(\theta_0, \epsilon)$ by Assumption 5, then there must exist a local maximum so that $n^{-1} \sum_{i=1}^n \frac{\partial \log p(X_i | \theta)}{\partial \theta} = 0$. Hence, for any small $\epsilon > 0$, there exists a sequence $\hat{\theta}_n(\epsilon)$ such that $P\left(\left|\hat{\theta}_n(\epsilon) - \theta_0\right| < \epsilon\right) \rightarrow 1$. To eliminate the dependence of the estimator on ϵ , we simply choose $\hat{\theta}_n$ that is closest to θ_0 for each n . The resulting sequence $\hat{\theta}_n^*$ satisfies $P\left(\left|\hat{\theta}_n^* - \theta_0\right| < \epsilon\right) \rightarrow 1$. \square

Example 6. Consider a random sample X_1, \dots, X_n from $N(\theta, 1)$ with density

$$p(x | \theta) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(x - \theta)^2}{2}\right\}, \quad x > 0, \theta \in \mathbb{R}.$$

The MLE is $\hat{\theta} = \bar{X}$.

1. We can show directly from the definition of consistency that

$$P(|\bar{X} - \theta| < \epsilon) = P\left(-\sqrt{n}\epsilon < \frac{\bar{X} - \theta}{\sqrt{1/n}} < \sqrt{n}\epsilon\right) \rightarrow 1.$$

2. We can also verify the conditions of the theorem. For example, for the identification assumption, if we have $\theta_1 \neq \theta_2$ such that $p(x | \theta_1) = p(x | \theta_2)$ for almost all x . Then, we must have

$$-\frac{(x - \theta_1)^2}{2} = -\frac{(x - \theta_2)^2}{2},$$

which means that θ_1 must be the solution of

$$\theta_1^2 - 2x\theta_1 - \theta_2^2 + 2x\theta_2 = 0.$$

If $x = 0$, we must have $\theta_1^2 = \theta_2^2$. But when $\theta_1 = -\theta_2$, we have

$$\theta_2^2 - 2x(-\theta_2) - \theta_2^2 + 2x\theta_2 = 4x\theta_2,$$

which cannot be zero unless $x = 0$. Thus, we must have $\theta_1 = \theta_2$.

Example 7 (Inconsistency of MLE: Ferguson's Example). Let X_1, \dots, X_n be with probability θ i.i.d. Uniform $(-1, 1)$, and be with probability $1 - \theta$ i.i.d. with a triangular distribution with pdf

$$\frac{1}{c(\theta)} \left(1 - \frac{|x - \theta|}{c(\theta)}\right), \text{ for } |x - \theta| \leq c(\theta),$$

where $c(\theta)$ is a continuous and decreasing function in θ with $c(0) = 1$ and $0 < c(\theta) \leq 1 - \theta$ for $0 < \theta < 1$. The parameter space is a compact set $\Theta = [0, 1]$. If $c(\theta) \rightarrow 0$ sufficiently fast as $\theta \rightarrow 1$, then $\hat{\theta}_n \xrightarrow{a.s.} 1$ whatever be the true value of $\theta \in \Theta$.

- In this example, there is no common support, since the triangular distribution part depends on the parameter.

- The inconsistency arises because the likelihood explodes if $c(\theta)$ is too small. Only one observation is enough to make it explode.

Theorem 4 (Strong Consistency of MLE). *Suppose that X_1, \dots, X_n are i.i.d., and satisfy*

1. *The parameter space Θ is a [compact set](#),*
2. *The density $p(x | \theta)$ with respect to μ is continuous in θ for all x ,*
3. *There exists a function $K(x)$ such that $E[|K(X)| | \theta_0] < \infty$ and*

$$U(x, \theta) = \log p(x | \theta) - \log p(x | \theta_0) \leq K(x),$$

for all x and θ ,

4. *for all $\theta \in \Theta$ and sufficiently small $\rho > 0$, $\sup_{|\theta' - \theta| < \rho} p(x | \theta')$ is measurable in x ,*
5. *$p(x | \theta) = p(x | \theta_0)$ almost everywhere with respect to μ implies $\theta = \theta_0$.*

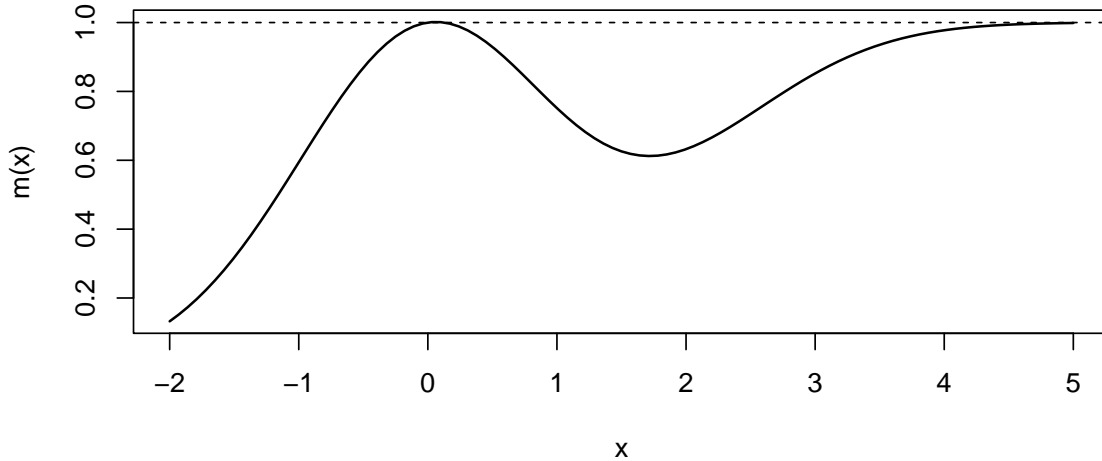
Then, any sequence of maximum likelihood estimator $\hat{\theta}_n$ of θ satisfies

$$\hat{\theta}_n \xrightarrow{a.s.} \theta.$$

In a general case, let $\hat{\theta}_n$ be an M-estimator that maximizes $M_n(\theta)$. We want our estimator to be [consistent](#): $d(\hat{\theta}_n, \theta_0) \xrightarrow{P} 0$, where $d(\cdot, \cdot)$ is a distance function. Heuristically,

1. [Pointwise convergence](#): for every θ , $M_n(\theta) \xrightarrow{P} M(\theta)$, where $M_n(\theta)$ is a random function and $M(\theta)$ is a deterministic function,
2. $\hat{\theta}_n = \arg \sup_{\theta} M_n(\theta)$,
3. $\theta_0 = \arg \sup_{\theta} M(\theta)$.

Then it is reasonable to expect $\hat{\theta}_n \xrightarrow{P} \theta_0$. However, pointwise convergence is too weak. An illustration is as follows.



Theorem 5 (Consistency of General M-Estimator). *Let M_n be random functions and let M be a fixed function of θ such that for every $\epsilon > 0$,*

$$\begin{aligned} \sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| &\xrightarrow{P} 0, \quad (\text{uniform convergence}) \\ \sup_{\theta: d(\theta, \theta_0) \geq \epsilon} M(\theta) &< M(\theta_0). \quad (\text{well-separated function}) \end{aligned}$$

Then, any sequence of estimators $\hat{\theta}_n$ with

$$M_n(\hat{\theta}_n) \geq M_n(\theta_0) - o_P(1) \quad (\text{nearly maximizer})$$

converges in probability to θ_0 .

Proof. The uniform convergence assumption implies $M_n(\theta_0) \xrightarrow{P} M(\theta_0)$, or equivalently

$$M(\theta_0) - o_P(1) \leq M_n(\theta_0) \leq M(\theta_0) + o_P(1).$$

Hence, by the nearly maximizer assumption,

$$M_n(\hat{\theta}_n) \geq M_n(\theta_0) - o_P(1) \geq M(\theta_0) - o_P(1).$$

This implies that

$$\begin{aligned} M(\theta_0) - M(\hat{\theta}_n) &\leq M_n(\hat{\theta}_n) + o_P(1) - M(\hat{\theta}_n) \\ &\leq \sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| + o_P(1) \xrightarrow{P} 0, \end{aligned}$$

where the convergence holds by the uniform convergence assumption.

The well-separated function assumption implies that, there exists $\eta(\epsilon) > 0$ such that for any $\theta \in \{\theta : d(\theta, \theta_0) \geq \epsilon\}$, we have $M(\theta) < M(\theta_0) - \eta$. If $\hat{\theta}_n \in \{\theta : d(\theta, \theta_0) \geq \epsilon\}$, then we have $M(\hat{\theta}_n) < M(\theta_0) - \eta$ and

$$P(d(\hat{\theta}_n, \theta_0) \geq \epsilon \mid \theta_0) \leq P(M(\hat{\theta}_n) < M(\theta_0) - \eta \mid \theta_0) = P(M(\theta_0) - M(\hat{\theta}_n) > \eta \mid \theta_0) \rightarrow 0,$$

where the convergence holds since we have shown above that $M(\theta_0) - M(\hat{\theta}_n) \xrightarrow{P} 0$. \square

Remark 3. If θ_0 is a unique maximizer of $M(\theta)$, then the well-separated function assumption means that the supremum is only attained at θ_0 . If $M_n(\hat{\theta}_n) \geq \sup_{\theta} M_n(\theta) - o_P(1)$, then $\hat{\theta}_n$ nearly maximizes M_n . But $M_n(\hat{\theta}_n) \geq M_n(\theta_0) - o_P(1)$ is enough for consistency.

Theorem 6 (Consistency of General Z-Estimator). *Let Ψ_n be random functions and let Ψ be a fixed function of θ such that for every $\epsilon > 0$,*

$$\begin{aligned} \sup_{\theta \in \Theta} \|\Psi_n(\theta) - \Psi(\theta)\| &\xrightarrow{P} 0, \quad (\text{uniform convergence}) \\ \inf_{\theta: d(\theta, \theta_0) \geq \epsilon} \|\Psi(\theta)\| &> 0 = \|\Psi(\theta_0)\|. \quad (\text{well separated function}) \end{aligned}$$

Then, any sequence of estimators $\hat{\theta}_n$ with

$$\Psi_n(\hat{\theta}_n) = o_P(1) \quad (\text{nearly a zero})$$

converges in probability to θ_0 .

We often do not need so strong assumptions as in the general case.

Lemma 3 (Consistency of Z-Estimator: Weaker Assumption). *Let Θ be a subset of the real line and let Ψ_n be random functions and let Ψ be a fixed function of θ such that*

$$\Psi_n(\theta) \xrightarrow{P} \Psi(\theta) \text{ for every } \theta. \quad (\text{pointwise convergence})$$

Assume that each map $\theta \mapsto \Psi_n(\theta)$ is continuous and has exactly one zero $\hat{\theta}_n$, or is nondecreasing with $\Psi_n(\hat{\theta}_n) = o_P(1)$. Let θ_0 be a point such that $\Psi(\theta_0 - \epsilon) < 0 < \Psi(\theta_0 + \epsilon)$ for every $\epsilon > 0$. Then, $\hat{\theta}_n \xrightarrow{P} \theta_0$.

Proof. There are two sets of assumptions. We prove them separately.

1. Because $\Psi_n(\theta)$ is a continuous function with a unique zero at $\hat{\theta}_n$, then if we know $\Psi_n(\theta_0 - \epsilon) < 0 < \Psi_n(\theta_0 + \epsilon)$, the solution must be between $\theta_0 - \epsilon$ and $\theta_0 + \epsilon$. This meant that

$$P(\Psi_n(\theta_0 - \epsilon) < 0 < \Psi_n(\theta_0 + \epsilon)) = P(\theta_0 - \epsilon < \hat{\theta}_n < \theta_0 + \epsilon).$$

By the assumption of pointwise convergence, we have $\Psi_n(\theta_0 - \epsilon) \xrightarrow{P} \Psi(\theta_0 - \epsilon)$ and $\Psi_n(\theta_0 + \epsilon) \xrightarrow{P} \Psi(\theta_0 + \epsilon)$. Thus, together with the assumption $\Psi(\theta_0 - \epsilon) < 0 < \Psi(\theta_0 + \epsilon)$, we further get

$$P(\Psi_n(\theta_0 - \epsilon) < 0 < \Psi_n(\theta_0 + \epsilon)) \rightarrow 1.$$

Hence,

$$P(\theta_0 - \epsilon < \hat{\theta}_n < \theta_0 + \epsilon) \rightarrow 1.$$

2. If we indeed have $\Psi_n(\hat{\theta}_n) = 0$, then a nondecreasing $\Psi_n(\theta)$ means that $\Psi_n(\theta_0 - \epsilon) \leq 0 \leq \Psi_n(\theta_0 + \epsilon)$ implies $\theta_0 - \epsilon \leq \hat{\theta}_n \leq \theta_0 + \epsilon$. We can use the same reasoning as in the first assumption set. Next we consider $0 \neq \Psi_n(\hat{\theta}_n) = o_P(1)$. For a $\eta > 0$, the nondecreasing $\Psi_n(\theta)$ means that

$$\begin{array}{l} \Psi_n(\theta_0 - \epsilon) < -\eta \\ \hat{\theta}_n \leq \theta_0 - \epsilon \end{array} \quad \text{implies} \quad \Psi_n(\hat{\theta}_n) \leq \Psi_n(\theta_0 - \epsilon) < -\eta.$$

Since $\Psi_n(\hat{\theta}_n) = o_P(1)$, then $P(\Psi_n(\hat{\theta}_n) < -\eta) \rightarrow 0$. If $\Psi_n(\theta_0 - \epsilon) < -\eta$ and $\Psi_n(\theta_0 + \epsilon) > \eta$, then $\hat{\theta}_n$ must satisfy $\theta_0 - \epsilon < \hat{\theta}_n < \theta_0 + \epsilon$ since $\Psi_n(\theta)$ is nondecreasing. Thus,

$$P(\Psi_n(\theta_0 - \epsilon) < -\eta \text{ and } \Psi_n(\theta_0 + \epsilon) > \eta) \leq P(\theta_0 - \epsilon < \hat{\theta}_n < \theta_0 + \epsilon).$$

By the assumption $\Psi(\theta_0 - \epsilon) < 0 < \Psi(\theta_0 + \epsilon)$ and pointwise convergence, we have

$$P(\Psi_n(\theta_0 - \epsilon) < -\eta \text{ and } \Psi_n(\theta_0 + \epsilon) > \eta) \rightarrow 1$$

for sufficiently small η .

□

Example 8 (Consistency of Median). The sample median of continuous random variable is a zero of the map

$$\theta \mapsto \Psi_n(\theta) = -\frac{1}{n} \sum_{i=1}^n \text{sign}(X_i - \theta).$$

- By LLN, we have pointwise convergence:

$$-\frac{1}{n} \sum_{i=1}^n \text{sign}(X_i - \theta) \xrightarrow{P} \Psi(\theta) = -E[\text{sign}(X - \theta)] = P(X < \theta) - P(X > \theta).$$

- $\Psi_n(\theta)$ is a nondecreasing function and $\Psi_n(\hat{\theta}_n) = 0$.
- The last assumption becomes

$$\begin{aligned} P(X < \theta_0 - \epsilon) - P(X > \theta_0 - \epsilon) &< 0 < P(X < \theta_0 + \epsilon) - P(X > \theta_0 + \epsilon), \\ \text{or } 2P(X < \theta_0 - \epsilon) + P(X = \theta_0 - \epsilon) &< 1 < 2P(X < \theta_0 + \epsilon) + P(X = \theta_0 + \epsilon). \end{aligned}$$

It holds for example if we have a continuous random variable and the population median is unique, i.e., $P(X < \theta_0 - \epsilon) < 0.5 < P(X < \theta_0 + \epsilon)$.

7.3 Asymptotic Normality

Theorem 7 (Cramér-Rao Conditions for MLE: Univariate). *Suppose that X_1, \dots, X_n are i.i.d., and satisfy*

1. *The parameter space Θ is an open set such that the true parameter value θ_0 is an interior point,*
2. *The distributions P_θ have common support $A = \{x : p(x | \theta) > 0\}$,*
3. *For every $x \in A$, the density $p(x | \theta)$ is three times differentiable with respect to θ , and the third derivative is continuous in θ ,*
4. *The integral $\int p(x | \theta) d\mu(x)$ can be twice differentiable under the integral sign,*
5. *The Fisher information $\mathcal{I}(\theta)$ satisfies $0 < \mathcal{I}(\theta) < \infty$,*
6. *There exists a function $M(x)$ such that*

$$\left| \frac{\partial^3 \log p(x | \theta)}{\partial \theta^3} \right| \leq M(x)$$

for all $x \in A$ and θ in a neighborhood of θ_0 , and that $E[M(X) | \theta_0] < \infty$.

Then, any consistent sequence $\hat{\theta}_n$ of roots of the likelihood equation satisfies

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N(0, \mathcal{I}^{-1}(\theta)).$$

Proof. Let $\ell(\theta) = \sum_{i=1}^n \log p(x_i | \theta)$ be the log-likelihood. By the third assumption, we get the Taylor's theorem:

$$0 = \frac{d\ell(\hat{\theta}_n)}{d\theta} = \frac{d\ell(\theta_0)}{d\theta} + (\hat{\theta}_n - \theta_0) \frac{d^2\ell(\theta_0)}{d\theta^2} + \frac{1}{2} (\hat{\theta}_n - \theta_0)^2 \frac{d^3\ell(\theta_n^*)}{d\theta^3},$$

where θ_n^* lies between θ_0 and $\hat{\theta}_n$. Thus,

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = \frac{n^{-1/2} \frac{d\ell(\theta_0)}{d\theta}}{-\frac{1}{n} \frac{d^2\ell(\theta_0)}{d\theta^2} - \frac{1}{2n} (\hat{\theta}_n - \theta_0) \frac{d^3\ell(\theta_n^*)}{d\theta^3}},$$

provided that the denominator is nonzero.

1. The numerator satisfies

$$n^{-1/2} \frac{d\ell(\theta_0)}{d\theta} = \sqrt{n} \left[\frac{1}{n} \sum_{i=1}^n \frac{d \log p(x_i | \theta_0)}{d\theta} - 0 \right] \xrightarrow{d} N(0, \mathcal{I}(\theta_0)),$$

by CLT since

$$\mathbb{E} \left[\frac{d \log p(x_i | \theta_0)}{d\theta} \mid \theta_0 \right] = 0, \quad \text{Var} \left[\frac{d \log p(x_i | \theta_0)}{d\theta} \mid \theta_0 \right] = \mathcal{I}(\theta_0)$$

by Assumptions 1, 2, 4, 5.

2. Regarding the first term in denominator,

$$\frac{1}{n} \frac{d^2 \ell(\theta_0)}{d\theta^2} = \frac{1}{n} \sum_{i=1}^n \frac{d^2 \log p(x_i | \theta_0)}{d\theta^2} \xrightarrow{P} \mathbb{E} \left[\frac{d^2 \log p(x_i | \theta_0)}{d\theta^2} \mid \theta_0 \right] = -\mathcal{I}(\theta_0),$$

by LLN and Assumptions 1, 2, 4, 5.

3. Regarding the second term in denominator,

$$\left| \frac{1}{n} \frac{d^3 \ell(\theta_n^*)}{d\theta^3} \right| = \left| \frac{1}{n} \sum_{i=1}^n \frac{d^3 \log p(x_i | \theta_n^*)}{d\theta^3} \right| \leq \frac{1}{n} \sum_{i=1}^n \left| \frac{d^3 \log p(x_i | \theta_n^*)}{d\theta^3} \right| \leq \frac{1}{n} \sum_{i=1}^n M(x_i),$$

by Assumption 6, if θ_n^* is in a small neighborhood of θ_0 . Our assumptions 1, 2, 3 implies that $\hat{\theta}_n \xrightarrow{P} \theta_0$ by Theorem 3. Hence, the probability that θ_n^* is in a small neighborhood of θ_0 approaches to 1. Note that by assumption $\mathbb{E}[M(X) \mid \theta_0] < \infty$, then LLN implies that

$$\frac{1}{n} \sum_{i=1}^n M(x_i) \xrightarrow{P} \mathbb{E}[M(X) \mid \theta_0],$$

that is,

$$\mathbb{P} \left(\mathbb{E}[M(X) \mid \theta_0] - \epsilon \leq \frac{1}{n} \sum_{i=1}^n M(x_i) \leq \mathbb{E}[M(X) \mid \theta_0] + \epsilon \right) \rightarrow 1.$$

Hence,

$$\mathbb{P} \left(\left| \frac{1}{n} \frac{d^3 \ell(\theta_n^*)}{d\theta^3} \right| \leq \mathbb{E}[M(X) \mid \theta_0] + \epsilon \right) \geq \mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n M(x_i) \leq \mathbb{E}[M(X) \mid \theta_0] + \epsilon \right) \rightarrow 1,$$

in other words, $\frac{1}{n} \frac{d^3 \ell(\theta_n^*)}{d\theta^3}$ is bounded in probability. This means that

$$\frac{1}{n} (\hat{\theta}_n - \theta_0) \frac{d^3 \ell(\theta_n^*)}{d\theta^3} = o_P(1).$$

Thus, we have reached

$$\sqrt{n} (\hat{\theta}_n - \theta_0) = \frac{n^{-1/2} \frac{d\ell(\theta_0)}{d\theta}}{\mathcal{I}(\theta_0) + o_P(1)} = \frac{n^{-1/2} \frac{d\ell(\theta_0)}{d\theta}}{\mathcal{I}(\theta_0)} + o_P(1).$$

□

Example 9. Consider $p(x | \theta) = \theta^{-1} \exp \{-\theta^{-1}x\}$, $\theta > 0$. The second-order derivative is

$$\frac{\partial^2 p(x | \theta)}{\partial \theta^2} = \left(\frac{x^2}{\theta^5} - \frac{4x}{\theta^4} + \frac{2}{\theta^3} \right) \exp \left\{ -\frac{x}{\theta} \right\},$$

and $\int \frac{\partial^2 p(x | \theta)}{\partial \theta^2} dx = 0$. In fact, this distribution belongs to the exponential family. Hence, we can change the order of integration and differentiation. Note that

$$\frac{\partial \log p(x | \theta)}{\partial \theta} = -\frac{1}{\theta} + \frac{x}{\theta^2}, \quad \frac{\partial^2 \log p(x | \theta)}{\partial \theta^2} = \frac{1}{\theta^2} - \frac{2x}{\theta^3}.$$

Then, the Fisher information is

$$\text{var} \left(\frac{\partial \log p(x | \theta)}{\partial \theta} \right) = \frac{1}{\theta^2} \in (0, \infty), \quad \text{or} \quad -E \left(\frac{\partial^2 \log p(x | \theta)}{\partial \theta^2} \right) = \frac{1}{\theta^2}.$$

The theorem above obtains the MLE by solving

$$\sum_{i=1}^n \frac{d \log p(x_i | \theta)}{d\theta} = 0.$$

Hence, it can be viewed as a Z-estimator. The following theorem considers the general Z-estimator under classic conditions.

Theorem 8 (Normality of Z-Estimator: Classic Condition). *Suppose that X_1, \dots, X_n are i.i.d., and consider $\Psi_n(\theta) = n^{-1} \sum_{i=1}^n \psi_\theta(X_i)$. Suppose that*

1. *For each θ in an open subset of Euclidean space, let $\theta \mapsto \psi_\theta(x)$ be twice continuously differentiable for every x .*
2. *Suppose that $E[\psi_{\theta_0}(X_1) | \theta_0] = 0$, $E[\|\psi_{\theta_0}(X_1)\|^2 | \theta_0] < \infty$ and that the matrix $E\left[\frac{d\psi_{\theta_0}(X_1)}{d\theta} \mid \theta_0\right]$ exists and is nonsingular.*
3. *Assume that the second-order partial derivatives are dominated by a fixed integrable function $\phi(x)$ for every θ in a neighborhood of θ_0 .*

Then every consistent estimator sequence $\hat{\theta}_n$ such that $\Psi_n(\hat{\theta}_n) = 0$, the sequence $\sqrt{n}(\hat{\theta}_n - \theta_0)$ is asymptotically normal with mean zero and covariance matrix

$$\left(E \left[\frac{d\psi_{\theta_0}(X_1)}{d\theta^T} \mid \theta_0 \right] \right)^{-1} E[\psi_{\theta_0}(X_1) \psi_{\theta_0}^\top(X_1) | \theta_0] \left\{ \left(E \left[\frac{d\psi_{\theta_0}(X_1)}{d\theta \dot{A} T} \mid \theta_0 \right] \right)^{-1} \right\}^\top.$$

7.4 Asymptotic Test

Suppose that

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, \sigma^2).$$

Then, we know

$$P \left(\hat{\theta} - \frac{\hat{\sigma}t}{\sqrt{n}} \leq \theta \leq \hat{\theta} + \frac{\hat{\sigma}t}{\sqrt{n}} \right) = P \left(-t \leq \frac{\sqrt{n}(\hat{\theta} - \theta)}{\hat{\sigma}} \leq t \right) \rightarrow P(-t \leq N(0, 1) \leq t).$$

Thus, we can use it as an asymptotic confidence interval.

If we want to test $H_0 : \theta = \theta_0$, we can construct a [Wald test](#). The Wald test statistic for testing $H_0 : \theta = \theta_0$ versus $H_1 : \theta \neq \theta_0$ is

$$Z_W = \frac{\sqrt{n}(\hat{\theta} - \theta)}{\hat{\sigma}}.$$

We reject H_0 if $|Z_W| > \lambda_{1-\alpha/2}$.

Definition 4 (LRT). The likelihood ratio test (LRT) statistic for testing $H_0 : \theta \in \Theta_0$ versus $H_1 : \theta \in \Theta_1$ is

$$\lambda(x) = \frac{\sup_{\theta \in \Theta_0} L(\theta | x)}{\sup_{\theta \in \Theta} L(\theta | x)} = \frac{L(\hat{\theta}_0 | x)}{L(\hat{\theta} | x)},$$

where $\hat{\theta}_0$ is the MLE when $\theta \in \Theta_0$, and $\hat{\theta}$ is the MLE when $\theta \in \Theta$. An LRT rejects H_0 if $\lambda(x)$ is too small.

Theorem 9. Suppose that X_1, \dots, X_n are iid following P_θ , where $\theta \in \Theta$ and Θ is an open set in \mathbb{R}^p . Suppose that,

1. $\frac{\partial \log p(x|\theta)}{\partial \theta}$ is well defined,
2. the true value θ_0 is the solution to $E \left[\frac{\partial \log p(x|\theta)}{\partial \theta} | \theta \right] = 0$,
3. the MLE is the solution to $\frac{\partial \log L(\theta; x)}{\partial \theta} = 0$ and is a consistent estimator of θ_0 ,
4. $E \left[\left\| \frac{\partial \log p(x|\theta)}{\partial \theta} \right\|^2 | \theta_0 \right] < \infty$, where $\| \cdot \|$ is the Euclidean norm,
5. $\frac{\partial^2 \log p(x|\theta)}{\partial \theta \partial \theta^T}$ is nonsingular and $E \left[\left\| \frac{\partial^2 \log p(x|\theta)}{\partial \theta \partial \theta^T} \right\| | \theta_0 \right] < \infty$,
6. the model is identified, i.e., $p(x | \theta_1) = p(x | \theta_2)$ almost everywhere under μ implies $\theta_1 = \theta_2$,
7. we can interchange the order of integration and differentiation such that $\text{Var} \left[\frac{\partial \log p(x|\theta)}{\partial \theta} | \theta \right] = -E \left[\frac{\partial^2 \log p(x|\theta)}{\partial \theta \partial \theta^T} | \theta \right]$.

Consider testing $H_0 : \theta \in \Theta_0$ versus $H_1 : \theta \in \Theta_1$. Suppose that the true value θ_0 satisfies H_0 . Then,

$$-2 \log \lambda(x) \xrightarrow{d} \chi_v^2,$$

where the degrees of freedom v is the difference in the number of free parameters.

Proof. We provide a heuristic proof for $H_0 : \theta = \theta_0$. Consider the Taylor expansion

$$\begin{aligned} \sum_{i=1}^n \log p(x_i | \theta_0) &\approx \sum_{i=1}^n \log p(x_i | \hat{\theta}) + \frac{\partial \sum_{i=1}^n \log p(x_i | \hat{\theta})}{\partial \theta} (\hat{\theta} - \theta_0) + \frac{1}{2} (\hat{\theta} - \theta_0)^T \frac{\partial^2 \sum_{i=1}^n \log p(x_i | \hat{\theta})}{\partial \theta \partial \theta^T} (\hat{\theta} - \theta_0) \\ &= \sum_{i=1}^n \log p(x_i | \hat{\theta}) + \frac{1}{2} (\hat{\theta} - \theta_0)^T \frac{\partial^2 \sum_{i=1}^n \log p(x_i | \hat{\theta})}{\partial \theta \partial \theta^T} (\hat{\theta} - \theta_0). \end{aligned}$$

Hence,

$$\begin{aligned} -2 \left\{ \sum_{i=1}^n \log p(x_i | \theta_0) - \sum_{i=1}^n \log p(x_i | \hat{\theta}) \right\} &\approx \sqrt{n} (\hat{\theta} - \theta_0)^T \left[-\frac{1}{n} \frac{\partial^2 \sum_{i=1}^n \log p(x_i | \hat{\theta})}{\partial \theta \partial \theta^T} \right] \sqrt{n} (\hat{\theta} - \theta_0) \\ &\approx N(0, \mathcal{I}^{-1}) \times \mathcal{I} \times N(0, \mathcal{I}^{-1}). \end{aligned}$$

□