جمهورية العراق
وزارة التعليم العالي والبحث العلمي
جامعة بابل/ كلية التربية

# حول الخواص التقاربية لمقدرات الأمكان الأعظم

بحث

مقدم الى قسم الرياضياتـ كلية التربيةـ جامعة بابل وهو جزء من متطلبات
نيل درجة الماجستير في علوم الرياضيات

من قبل

## رواسي عدنان حميد المسعودي

بإشراف

## الاستاذ المساعد الدكتور سعد عبد ماضي

1424هـ        2003 م

Republic of Iraq,
Ministry of Higher Education,
and Scientific Research,
Babylon University,
College of Education

# ON ASYMPTOTIC PROPERTIES OF MAXIMUM LIKELIHOOD ESTIMATORS

A Research
Submitted to the Department of Mathematics- College of Education-
University of Babylon in Partial fulfillment of the requirements
for the Degree of Master of Science
in Mathematics

By
Rawassy A. Hameed Al-Masa'oodi

Supervised by
Assist. Prof. Dr. Sa'ad A. Madhi

2003

# ABSTRACT

A central role is played in large sample parametric inference by the maximum likelihood estimation of the parameter of interest. This project deals with aspects of large sample behavior of the maximum likelihood estimators such as consistency and asymptotic normality.

For rigorous and complete proofs, analytic details and regularity conditions are introduced.

Illustrative examples and some concluding remarks are also given.

# حول الخواص التقاربية لمقدرات الأمكان الأعظم

# الخلاصة

يلعب مقدر الامكان الاعظم للمعلمة مثار الاهتمام دورا مركزياً في الاستدلال المعلمي في حالة العينة الكبيرة.

هذا البحث يتناول مظاهر سلوك مقدر الامكان الاعظم في حالة العينة الكبيرة، مثل الاتساق وسلوك التوزيع الطبيعي تقريباً وغيرها.

ولكي تكون البراهين كاملة ودقيقة تم عرض التفاصيل التحليلية المطلوبة مع الافتراضات المناسبة لها.

وتضمن البحث ايضاً بعض الامثلة التوضيحية للنتائج بالاضافة الى بعض الملاحظات حول هذه النتائج.

# ACKNOWLEDGEMENTS

# الاهداء

اهدي بحثي هذا الى من شجعني دائماً واحاطني بالعطف والحنان ابداً وقام بمساعدتي قبل واثناء اتمامي لبحثي هذا وهو عائلتي (ابي وامي واخوتي واختي) شيء بسيط لما قاموا به لاجلي.

وهو ايضاً مهداة الى اقربائي عائلة السيد عبود السيد مهدي الشلاه بجميع افرادها لدعمهم اللامحدود لي ومساعدتهم المستمرة لي.

ومن الله التوفيق

# CONTENTS

# <u>CERTIFICATION</u>

We certify that we have read this research entitled " On Asymptotic Properties of Maximum Likelihood Estimator" and, as an Examining Committee, we examined the student in the contents, and what is connected with, and that in our opinion it meets the standard of a research for the degree of Master of Science in Mathematics.

**Chairman**

**Signiture:**

**Name:**

**Date:**

**Member**                                    **Member**

**Signiture:**                                 **Signiture:**

**Name:**                                      **Name:**

**Date:**                                      **Date:**

**Member**     **(Supervisor)**

**Signiture:**

**Name:**

**Date:**

Approved by the Dean of the College of Education

**Signature:**

**Name:**

**Date:**

بسم الله الرحمن الرحيم

وقل اعملوا فسيرى الله عملكم

ورسوله والمؤمنون

صدق الله العظيم

سورة التوبة
الآية (105)

# INTRODUCTION

The chief justification for the method of maximum likelihood estimation lies in its near-optimal properties for large samples.

This project deals with the asymptotic properties of the maximum likelihood estimators such as consistency and asymptotic normality. To do this rigorously requires analytic details and suitable regularity conditions. Some examples are presented to illustrate the results. Concluding remarks are also given.

A summary of the chapters is as follows:

Chapter two provides some concepts and results that are needed in later chapters.

The large –sample properties of the maximum likelihood estimators are studied in chapter three.

Chapter four contains examples illustrating the results in chapter three.

Finally, in chapter five some comments on the large–sample properties of the maximum likelihood estimator are presented.

# PRELIMINARIES

In this chapter, we cover thoroughly many concepts and results which find use in later chapters. We refer back to this chapter if and as the need arise.

## 1.1- Definitions and results from Analysis and Topology

## 1.1.1- Limit of sequences     [1]

Let $\{A_n\}$ be a sequence of sets. The set of all points $\omega$ in the set $\Omega$ that belong to $A_n$ for infinitely many values of $n$ is called the limit superior of the sequence $\{A_n\}$ and is written as $\limsup\limits_{n\to\infty} A_n$.

The set of all points that belong to $A_n$ for all but a finite number of values of $n$ is called the limit inferior of the sequence $\{A_n\}$ and is written as $\liminf\limits_{n\to\infty} A_n$.

If $\limsup\limits_{n\to\infty} A_n = \liminf\limits_{n\to\infty} A_n$,
we say that limit exists of the sequence $\{A_n\}$ and we denote it by $\lim\limits_{n\to\infty} A_n$.

For an increasing sequence $A_1 \subset A_2 \subset ...,$ $\lim\limits_{n\to\infty} A_n = \bigcup\limits_{n} A_n$.

For a decreasing sequence $A_1 \supset A_2 \supset ...,$ $\lim\limits_{n\to\infty} A_n = \bigcap\limits_{n} A_n$.

In the case of an arbitrary sequence of sets $A_1, A_2,...$ , we have $\liminf\limits_{n\to\infty} A_n = \bigcup\limits_{n=1}^{\infty} \bigcap\limits_{k=n}^{\infty} A_k,$ $\limsup\limits_{n\to\infty} A_n = \bigcap\limits_{n=1}^{\infty} \bigcup\limits_{k=n}^{\infty} A_k.$

## 1.1.2- Upper Semi Continuous Function     [4]

A real–valued function $f(\theta)$ defined on the set $\Theta$ is said to be upper semi continuous on $\Theta$, if for all $\underline{\theta}$ in $\Theta$ and for any sequence $\underline{\theta}_n$ in $\Theta$, such that $\underline{\theta}_n \to \underline{\theta}$ , we have

$$\lim_{n\to\infty} \sup f(\underline{\theta}_n) \le f(\underline{\theta})$$

or, equivalently, if for all $\underline{\theta}$ in $\Theta$,

$$\sup_{|\underline{\theta}'-\underline{\theta}|<\rho} f(\underline{\theta}') \to f(\underline{\theta}) \quad \text{as} \quad \rho \to 0 .$$

## 1.1.3- Taylor Expansion    [4]

If $f : R^k \to R$ and if $f''(x)$ is continuous in the sphere $\{\underline{x} : |\underline{x} - \underline{x}_\circ| < r\}$, then for $|\underline{t}| < r$,

$$f(\underline{x}_\circ + \underline{t}) = f(\underline{x}_\circ) + f'(\underline{x}_\circ)\underline{t} + \underline{t}^T \int_0^1\int_0^1 \lambda\, f''(\underline{x}_\circ + u\lambda\,\underline{t})\,du\,d\lambda\,\underline{t}$$

where $R^k$ is the k-dimensional Euclidean space.

## 1.1.4- Convex Function    [3]

Let $U \subseteq R^d$ and $\underline{x}_1$ and $\underline{x}_2 \in R^d$, then the point $\underline{x} = (1-\lambda)\underline{x}_1 + \lambda\,\underline{x}_2,\ 0 \le \lambda \le 1$, is called the convex linear combination of $\underline{x}_1$ and $\underline{x}_2$. It is any point lying on the line segment joining $\underline{x}_1$ and $\underline{x}_2$. A set $U \subseteq R^d$ is said to be convex, if the convex linear combination of any two points in $U$ belongs to $U$.

In other words, $U$ is convex if $\underline{x}_1, \underline{x}_2 \in U$, implies that $\underline{x} \in U$, where $\underline{x} = (1-\lambda)\underline{x}_1 + \lambda\,\underline{x}_2,\ 0 \le \lambda \le 1$.

If $\underline{x} \in U \subset R^d$, where $U$ is a convex set, then the function $f(\underline{x})$ is said to be convex if for any two points $\underline{x}_1$ and $\underline{x}_2$ in $U$,

$$f(\underline{x}) \le (1-\lambda)\,f(\underline{x}_1) + \lambda\,f(\underline{x}_2) \qquad , 0 \le \lambda \le 1$$

for every $\underline{x} = (1-\lambda)\underline{x}_1 + \lambda\,\underline{x}_2$.

The function is said to be concave if the inequality sign is reversed or if $-f(\underline{x})$ is convex. If the inequality is strict, then $f$ is called strictly convex function.

## 1.1.5- σ-algebra (σ-field)    [1]

A non-empty class $F$ of subsets of a set $\Omega$ which is closed under the formulation of countable unions and complements and contains $\emptyset$ is known as a σ-field.

## 1.1.6- Measure    [1]

A non- negative extended real-valued set function $\mu$ defined on a class $F$ of subsets of a set $\Omega$ is called a measure on $F$, if $\mu(A) \geq 0$, for all $A \in F$, and $F$ is $\sigma - \text{algebra}$, and

$$\mu \left( \bigcup_{n=1}^{\infty} A_n \right) = \sum_{n=1}^{\infty} \mu(A_n).$$

## 1.1.7- Measurable Space    [1]

A measurable space is a pair $(\Omega, F)$ when $\Omega$ is a set and $F$ is a σ-field of subsets of $\Omega$. Any member of $F$ is called a measurable set.

## 1.1.8- Measurable Function    [1]

Let $(\Omega, F)$ and $(\Omega', F')$ be two measurable spaces.

A function $f : \Omega \to \Omega'$ is said to be measurable (relative to $F$ and $F'$) if $f^{-1}(A) \in F$, for every $A \in F'$.

## 1.1.9- σ-Finite Measure    [1]

4

An extended real valued set function $v$ defined on a class $F$ of subset of a set $\Omega$ is said to be σ-finite measure if for each $A \in F$ , there is a sequence of sets $A_n \in F$ , such that

$$A \subset \bigcup_n A_n \quad , \text{ and } \quad v(A_n) \text{ is finite for all } n.$$

i.e $v(A_n) < \infty.$

## 1.1.10- Topological Space    [7]

Let $\Omega$ be a non-empty set. A class $F$ of subsets of $\Omega$ is a topology on $\Omega$ iff $F$ satisfies the following axioms:

1- $\Omega$ and Ø belong to $F$.

2- The union of any number of sets in $F$ belongs to $F$.

3- The intersection of any two sets in $F$ belongs to $F$. The members of $F$ are then called $F$-open sets, or simply open sets. The pair ($\Omega$, $F$) is called a topological space.

## 1.1.11- Open Cover    [7]

Let $C = \{G_\alpha : \alpha \in \wedge\}$ be open collection of subsets of a topological space X and a subset A of X, we say that C is an open cover of the space X if $X = \cup\{G_\alpha : \alpha \in \wedge\}$ .

## 1.1.12- Finite Sub-Cover    [7]

A subcollection $C'$ of $C$ , such that $C'$ covers A. An open cover C of A is said to be finite sub-cover if it consists of finite number of open sets $C'$ .v

## 1.1.13- Compact Set    [7]

A subset A of a topological space X is said to be compact iff every open cover of A has a finite sub-cover. More explicitly, that is iff for every collection $\{G_\alpha : \alpha \in \wedge\}$ of open sets for which $A \subset \cup \{G_\alpha : \alpha \in \wedge\}$, there exists finitely many sets $G_{\alpha_1}, G_{\alpha_2}, ..., G_{\alpha_n}$ among the $G_\alpha$'s such that:

$$A \subset \cup G_{\alpha_1} \cup ... \cup G_{\alpha_n}.$$

## 1.2- Definitions and results from Probability and Statistics
## 1.2.1-Modes of Convergence        [3], [4]

For a random vector $\underline{x} = (x_1, x_2, ..., x_d) \in R^d$, the distribution function of $\underline{x}$, defined for $\underline{x} = (x_1, x_2, ..., x_d) \in R^d$, is denoted by:

$$F_{\underline{x}}(x) = p(\underline{x} \leq \underline{x}) = p(x_1 \leq x_1, ..., x_d \leq x_d) \in R^d.$$

The Euclidean norm of $\underline{x} = (x_1, x_2, ..., x_d) \in R^d$ is denoted by

$$|\underline{x}| = (x_1^2 + x_2^2 + ... + x_d^2)^{\frac{1}{2}}.$$

Let $\underline{x}, \underline{x}_1, \underline{x}_2, ...$ be random vectors with values in $R^d$. we have the following modes of convergence.

*(a) Convergence in Distribution*

The sequence $\{x_n\}$ converges in law to $\underline{x}$, $x_n \overset{L}{\to} \underline{x}$, if $F_{\underline{x}_n}(x) \to F_{\underline{x}}(x)$ as $n \to \infty$, for all points $\underline{x}$ at which $F_{\underline{x}}(x)$ is continuous. It is sometimes called convergence in distribution or weak convergence.

*(b) Convergence in Probability*

The sequence $\{x_n\}$ converges in probability to $\underline{x}$, $x_n \xrightarrow{p} \underline{x}$ if for every $\varepsilon > 0$,

$$p\{|x_n - \underline{x}| > \varepsilon\} \to 0, \quad \text{as} \quad n \to \infty.$$

*(c) Convergence in $r^{th}$ mean*

For a real number $r > 0$, the sequence $\{x_n\}$ converges in the $r^{th}$ mean to $\underline{x}$, $x_n \xrightarrow{r} \underline{x}$, if

$$E|x_n - \underline{x}|^r \to 0, \quad \text{as} \quad n \to \infty.$$

where $E$ is the expectation value.

*(d) Convergence almost surely*

The sequence $\{x_n\}$ converges almost surely to $\underline{x}$, $x_n \xrightarrow{a.s} \underline{x}$, if $\quad p\{\lim_{n \to \infty} x_n = \underline{x}\} = 1$.

Almost sure converges is sometimes called convergence with probability 1 (w.p.1) or strong convergence.

Equivalently, $x_n \xrightarrow{a.s} x$ iff for every $\varepsilon > 0$,

$$p\{|x_m - \underline{x}| < \varepsilon, \text{for all} \quad m \geq n\} \to 1 \quad \text{as} \quad n \to \infty.$$

**1.2.2-Fatou-Lebesgue Theorem    [1], [4]**

If $x_n \xrightarrow{a.s} x$ and if for all $n$, $\underline{x}_n \geq \underline{y}$, for some random vector $\underline{y}$ with $E|\underline{y}| < \infty$, then $\quad \liminf_{n \to \infty} E(\underline{x}_n) \geq E(\underline{x})$.

**1.2.3- Montone Convergence Theorem    [4]**

If $\quad 0 \leq \underset{\sim}{x}_1 \leq \underset{\sim}{x}_2 \leq \ldots \quad$ and $\quad \underset{\sim}{x}_n \overset{a.s}{\to} \underset{\sim}{x}$ , then

$$E(\underset{\sim}{x}_n) \to E(\underset{\sim}{x}).$$

## 1.2.4- Lebesgue Dominated Convergence Theorem   [4]

If $\quad \underset{\sim}{x}_n \overset{a.s}{\to} \underset{\sim}{x} \quad$ and if $\quad |\underset{\sim}{x}_n| \leq \underset{\sim}{y}$, for some random vector $\underset{\sim}{y}$

with $E|\underset{\sim}{y}| < \infty$, then $E(\underset{\sim}{x}_n) \to E(\underset{\sim}{x})$, such that $E(\underset{\sim}{x})$ exists.

In these theorems, $\underset{\sim}{x}$ , $E(\underset{\sim}{x}_n)$ and $E(\underset{\sim}{x})$ may take the value $+\infty$.

## 1.2.5- The Law of Large Numbers    [6]

Let $\underset{\sim}{x}_1, \underset{\sim}{x}_2, \ldots$ be i.i.d (independent identically distributed) random vectors, and let

$$\overline{\underset{\sim}{x}}_n = \left(\frac{1}{n}\right) \sum_{j=1}^{n} \underset{\sim}{x}_j.$$

(a) If $\quad E|\underset{\sim}{x}| \leq \infty$ , then $\overline{\underset{\sim}{x}}_n \overset{p}{\to} \mu = E(\underset{\sim}{x})$. This is the weak law of large numbers (WLLN).

(b) $\overline{\underset{\sim}{x}}_n \overset{a.s}{\to} \mu$ iff $E|\underset{\sim}{x}| < \infty$ and $\mu = E(\underset{\sim}{x})$. This is the strong law of large number (SLLN).

(c) If $E|\underset{\sim}{x}|^2 < \infty$, then $\overline{\underset{\sim}{x}}_n \overset{qm}{\to} \mu = E(\underset{\sim}{x})$

($qm$ means in quadratic mean).

## 1.2.6- Central Limit Theorem   [5]

Let $\underset{\sim}{x}_1, \underset{\sim}{x}_2, \ldots,$ be i.i.d random vectors with mean $\underset{\sim}{\mu}$ and finite covariance matrix, $B$. Then $\quad \sqrt{n}(\overline{\underset{\sim}{x}}_n - \underset{\sim}{\mu}) \overset{L}{\to} N(\underset{\sim}{0}, B)$.

## 1.2.7- The Continuity Theorem    [4]

$$x_n \xrightarrow{L} x \quad \text{iff} \quad \varphi_{x_n}(t) \to \varphi_x(t), \quad \text{for all} \quad t \in R^d$$

where $\varphi_x$ is the characteristic function of $x$.

## 1.2.8- Jensen's Inequality  [3]

Let $x$ be a random vector with mean $E(x)$, and let $f(.)$ be convex function, then

$$E[f(x)] \geq f[E(x)].$$

## 1.2.9- Slutsky Theorems  [4]

(a) If $\{x_n\}$ be a sequence of random vectors, $x_n \in R^d$, $x_n \xrightarrow{L} x$

and if $f : R^d \to R^k$ is such that

$p\{x \in c(f)\} = 1$ , where $c(f)$ is the continuity set of $f$,

then $f(x_n) \xrightarrow{L} f(x)$.

(b) If $x_n \xrightarrow{L} x$ and $(x_n - y_n) \xrightarrow{p} 0$, then $y_n \xrightarrow{L} x$ .

(c) If $x_n \in R^d$, $y_n \in R^k$, $x_n \xrightarrow{L} x$ and $y_n \xrightarrow{L} c$,

then $\begin{pmatrix} x_n \\ y_n \end{pmatrix} \xrightarrow{L} \begin{pmatrix} x \\ c \end{pmatrix}$.

## 1.2.10- Consistency in Statistical Estimation  [2], [6]

In such problems the underlying probability, $P_{\theta_\circ}$, depends upon a parameter $\theta_\circ \in \Theta$ in $R^k$, where $\Theta$ is the parameter space.

Let $\{\tilde{\theta}_n\}$ be a sequence of random vectors, considered as estimators of $\theta_\circ$, the true value of the parameter.

We say that $\{\tilde{\theta}_n\}$ is weakly consistent sequence estimators

of $\theta_\circ$ if $\tilde{\theta}_n \xrightarrow{p} \theta_\circ$ for all $\theta_\circ \in \Theta$.

9

where $p = p_{\theta_\circ}$ is the "true" probability distribution. This is sometimes called consistency in probability. We may similarly define strong consistency ($\tilde{\theta}_n \overset{a.s}{\to} \theta_\circ$) or consistency in quadratic mean ($\tilde{\theta}_n \overset{qm}{\to} \theta_\circ$), both of which imply (weak consistency).

The weak (strong) law of large numbers states that the sample mean is a weakly (strongly) consistent estimator of the population mean. That is:

If $\underset{-}{x}_1, \underset{-}{x}_2, \ldots, \underset{-}{x}_n$ are i.i.d random vectors with common finite $\underset{-}{\mu}$ and $\overline{x}_n = \frac{1}{n} \sum_{i=1}^{n} \underset{-}{x}_i$,

then $\overline{x}_n \overset{p}{\to} \mu$ as $n \to \infty$ (weak consistency),

and $\overline{x}_n \overset{a.s}{\to} \mu$ as $n \to \infty$ (strong consistency).

## 1.2.11-Maximum likelihood Estimator     [3]

Let $x_1, x_2, \ldots, x_n$ be identically independent random sample with density $f(x/\theta)$ with respect to some σ-finite measure $v$ (usually Lebesgue measure or counting measure) where $\theta \in \Theta \subseteq R^k$. The function

$$L_n(\theta) = L_n(\theta, x) = \prod_{i=1}^{n} f(x_i, \theta) \qquad \ldots\ldots(1)$$

Considered as a function of $\theta$ is called the likelihood function, where the observed values of $x_1, x_2, \ldots, x_n$ are $x_1, x_2, \ldots, x_n$. The principle of maximum likelihood estimation consists of choosing as an estimator of $\theta$ any function $\hat{\theta}_n = \hat{\theta}_n (x_1, \ldots, x_n)$ such that

$$L_n(\hat{\theta}_n) = \sup_{\theta \in \Theta} L_n(\theta_n) \qquad \ldots\ldots(2)$$

If $\hat{\theta}$ satisfying (2) exists, we call it a maximum likelihood estimate (MLE).

It is convenient to work with logarithm of the likelihood function. $\ell_n(\theta) = \log L_n(\theta)$. Since log is monotone function.

$$\ell_n(\hat{\theta}_n) = \sup_{\theta \in \Theta} \ell_n(\theta)$$

When $(\frac{\partial}{\partial \theta})f(x,\theta)$ exists, one can seek the MLE $\hat{\theta}_n$ as

$$\frac{\partial L_n(\theta)}{\partial \theta} = 0 \quad \text{and} \quad \frac{\partial^2 L_n(\theta_n)}{\partial \theta^2} < 0$$

or equivalently

$$\frac{\partial \ell_n(\theta)}{\partial \theta} = 0 \quad \text{and} \quad \frac{\partial^2 \ell_n(\theta_n)}{\partial \theta^2} < 0 \ldots\ldots\ldots..(3)$$

Where $\quad \dfrac{\partial \ell_n(\theta)}{\partial \theta} = \dfrac{\partial \log L_n(\theta)}{\partial \theta} = \sum_1^n \dfrac{\partial}{\partial \theta} \log f(x_i,\theta) = 0 .$

Equation (3) is usually referred to as the likelihood equation.

## 1.2.12- Fisher Information        [8]

Fisher's information function plays an important role in the large sample theory of MLE.

Let

$\psi(x,\theta) = (\frac{\partial}{\partial \theta} \log f(x/\theta))^T$, a $k$ dimensional vector and

$\psi'(x,\theta) = \dfrac{\partial^2}{\partial \theta^2} \log f(x/\theta)$, a $k \times k$ matrix.

The Fisher information matrix for a single observation is defined as $B(\theta) = E_\theta(\psi(x,\theta)\psi(x,\theta)^T)$  a $k \times k$ matrix.

Assuming that the partial derivative with respect to $\theta$ can be passed under the integral sign in

$$\int f(x/\underline{\theta})dv(x) = 1 \quad ; \quad \text{we find}$$

$$E_{\underline{\theta}}(\psi(\underline{x},\underline{\theta})) = \int \frac{\partial/\partial\underline{\theta}\; f(x/\underline{\theta})}{f(x/\underline{\theta})} f(x/\underline{\theta})dv(x)$$

$$= \int \frac{\partial}{\partial\underline{\theta}} f(x/\underline{\theta})dv(x) = \underline{0}$$

so that $B(\underline{\theta})$ is in fact the covariance matrix of $\psi$,

$$B(\underline{\theta}) = \text{var}_{\underline{\theta}}(\psi(\underline{x},\underline{\theta}))$$

If the second partial derivatives with respect to $\underline{\theta}$ can be passed

under the intgeral sign, then

$$\int \left(\frac{\partial^2}{\partial\underline{\theta}^2}\right) f(x/\underline{\theta})dv(x) = \underline{0}, \text{ and}$$

$$E_{\underline{\theta}}(\psi'(\underline{x},\theta)) = \int \left[\frac{\partial}{\partial\underline{\theta}} \frac{(\partial/\partial\underline{\theta})f(x/\underline{\theta})}{f(x/\underline{\theta})}\right] f(x/\underline{\theta})dv(x)$$

$$= \int \frac{f(x/\underline{\theta})(\partial^2/\partial\underline{\theta}^2)f(x/\underline{\theta}) - ((\partial/\partial\underline{\theta})f(x/\underline{\theta}))^T(\partial/\partial\underline{\theta})f(x/\underline{\theta})}{f(x/\underline{\theta})^2} \times$$

$$f(x/\underline{\theta})dv(x)$$

$$= 0 - \int \psi(x,\underline{\theta})\psi(x,\underline{\theta})^T f(\underline{x}/\underline{\theta})dv(x) .$$

Thus

$$B(\underline{\theta}) = -E_{\underline{\theta}}(\psi'(\underline{x},\underline{\theta})) .$$

# ASYMPTOTIC PROPERTIES OF

# MAXIMUM LIKELIHOOD ESTIMATORS

The main justification of the method of maximum likelihood is a "large-sample" one, which shows that when an observation provides lots of information about an unknown parameter, the method utilizes essentially all of this information. In this chapter, we study the large-sample (asymptotic) properties of the MLE.

## 2.1- Consistency of MLE

A key property of the maximum likelihood estimators is their consistency under certain assumptions. The proof of consistency of the MLE that we will present in theorem (3.1.1) is mainly based on the ideas given in Wald (1949), Wolfowiz (1949). We start with the following lemma.

Let $f_0(x)$ and $f_1(x)$ be densities with respect to a $\sigma-$finite measure $v$.

The kullback- Leibler information number is defined as

$$K(f_0, f_1) = E_0\left( \log \frac{f_0(\mathrm{x})}{f_1(\mathrm{x})} \right) = \int (\log \frac{f_0(x)}{f_1(x)}) f_0(x) dv(x).$$

In this expression, $\log \dfrac{f_0(x)}{f_1(x)}$ is defined as $+\infty$ if $f_1(x) = 0$ and $f_0(x) > 0$ ,so the expectation could be $+\infty$. Although $\log \dfrac{f_0(x)}{f_1(x)}$ is defined as $-\infty$ when $f_1(x) > 0$ and $f_0(x) = 0$, the

integrand, $\log\left(\dfrac{f_0(x)}{f_1(x)}\right)f_0(x)$, is defined as zero in this case. $K(f_0, f_1)$ is a measure of the ability of the likelihood ratio to distinguish between $f_0$ and $f_1$ when $f_0$ is true.

## Lemma (Shannon-kolmogorov Information Inequality)

Let $f_0(x)$ and $f_1(x)$ be densities with respect to $v$. Then

$$K(f_0, f_1) = E_0\left(\log \frac{f_0(\mathrm{x})}{f_1(\mathrm{x})}\right) = \int \log \frac{f_0(x)}{f_1(x)} f_0(x)dv(x) \geq 0$$

With equality if and only if $f_1(x) = f_0(x)$ (a.e $dv$).

*Proof:*

Since $\log x$ is strictly convex. Jensen's inequality implies

$$-K(f_0, f_1) = E_0\left(\log \frac{f_1(\mathrm{x})}{f_0(\mathrm{x})}\right) \leq \log E_0 \frac{f_1(\mathrm{x})}{f_0(\mathrm{x})} \ ,\ldots\ldots\ldots (1)$$

With equality if and only if $\dfrac{f_1(\mathrm{x})}{f_0(\mathrm{x})}$ is a constant with probability 1 when x has density $f_0$. That is

if $\quad \dfrac{f_1(x)}{f_0(x)} = c$, then $\quad E(\log c) = \log c$.

But

$$E_0 \frac{f_1(\mathrm{x})}{f_0(\mathrm{x})} = \int \frac{f_1(x)}{f_0(x)} f_0(x)dv(x) = \int_{s_0} f_1(x)dv(x) \leq 1 ,\ldots\ldots (2)$$

Where $S_0 = \{x: f_0(x) > 0\}$, with equality iff $S_0$ has probability 1 under $f_1(x)$. The combination of (1) and (2) gives

$$-K(f_0, f_1) = \log E\left(\frac{f_1(\mathrm{x})}{f_0(\mathrm{x})}\right) \leq \log 1 = 0$$

or

$K(f_0, f_1) \geq 0$, With equality iff $f_1(x) = f_0(x)$ (a.e. $dv$).

## *Theorem (2.1.1)*

Let $x_1, x_2, \ldots, x_n$ be i.i.d with density $f(x/\theta)$, where $\underline{\theta} \in \Theta$, $\Theta$ is the parameter space.

Let $\theta_\circ$ denotes the true value of the parameter $\theta$.

$\{\hat{\theta}_n\}$ is the sequence of maximum likelihood estimators of the parameter $\theta$.

For the validity of the proof of this theorem. We assume the following assumptions called (regularity conditions):

1- $\Theta \subseteq \mathbb{R}^k$ is compact,

2- $f(x/\theta)$ is upper semi continuous in $\theta$ for all $x$,

3-there exists a function $K(x)$ such that $E_{\theta_\circ}|K(\mathrm{x})| < \infty$, and

$$W(x,\theta) = \log \frac{f(x/\theta)}{f(x/\theta_\circ)} \leq K(x) \text{ , for all } x \text{ and } \theta,$$

4-for all $\underline{\theta} \in \Theta$ and for all sufficiently small $\rho > 0$,

$$\sup_{|\theta'-\theta|<\rho} f(x/\theta') \text{ is measurable in } x,$$

5-(Identifiability), $f(x/\theta) = f(x/\theta_\circ)$ (a.e dv) $\Rightarrow \theta = \theta_\circ$,

if these conditions are satisfied, and for each $n \geq 1$ the MLE is unique, then, for any sequence of MLE $\hat{\theta}_n$ of $\theta$, $\hat{\theta}_n \xrightarrow{a.s} \theta_\circ$.

i.e. For any $\varepsilon > 0$ , and for all $\theta_\circ \in \Theta$

$$\lim_{n\to\infty} p\{|\hat{\theta}_n - \theta_\circ| < \varepsilon\} = 1 \quad \text{holds.}$$

i.e The MLE $\hat{\theta}_n$ is consistent.

## Proof:

Let $f(x,\theta,\rho) = \sup_{|\theta'-\theta|<\rho} W(x,\theta)$ .

Then $f(x,\theta_\circ,\rho)$ is measurable in $x$ for all sufficiently small $\rho > 0$ by (4).

15

By (3), $f(x,\theta,\rho)$ is bounded above by an integrable function $K(x)$, and $f(x,\theta_{\circ},\rho) \to W(x,\theta)$, as $\rho \to 0$ by (2) (since $W(x,\theta)$ is upper semi continuous in $\theta$ for all $x$.

Therefore, by the Monotone Convergence Theorem, we have

$$\int f(x,\theta,\rho)f(x)dv(x) \to \int W(x,\theta)f(x)dv(x) = \mu(\theta) \quad \text{as} \quad \rho \to 0.$$

Let $\varepsilon > 0$. Find $\rho_\theta$ for each $\theta$, so that

$$\int f(x,\theta,\rho_\theta)f(x)dv(x) < \mu(\theta) + \varepsilon.$$

The spheres $S(\theta,\rho_\theta) = \{\theta' : |\theta' - \theta| < \rho_\theta\}$ cover $\Theta$.

Since $\Theta$ is compact, then there exists a finite sub-cover, say $\Theta \subset \overset{m}{\underset{1}{\cup}} S(\theta_j, \rho_{\theta_j})$. For each $\underline{\theta} \in \Theta$ there exists an index j such that $\theta \in S(\theta_j, \rho_{\theta_j})$. From the definition of $f(x,\theta,\rho)$,

$$W(x,\theta) \leq f(x,\theta_j,\rho_{\theta_j}) \text{ for all } x. \text{ Hence,}$$

$$\frac{1}{n}\sum_{i=1}^{n} W(x_i,\theta) \leq \frac{1}{n}\sum_{1}^{n} f(x_i,\theta_j,\rho_{\theta_j}).$$

So that

$$\sup_{\underline{\theta}\in\Theta} \frac{1}{n}\sum_{1}^{n} W(x_i,\theta) \leq \sup_{1\leq j\leq m} \frac{1}{n}\sum_{1}^{n} f(x_i,\theta_j \rho_{\theta_j}).$$

Now, apply the SLLN to $\frac{1}{n}\sum_{1}^{n} f(x_i,\theta_j,\rho_{\theta_j})$,

$$p\{\lim_{n\to\infty} \frac{1}{n}\sum_{1}^{n} f(x_i,\theta_j,\rho_{\theta_j}) \leq \mu(\theta_j) + \varepsilon, \quad \text{for } j = 1,...,m\} = 1,$$

$$p\{\limsup_{n\to\infty} \sup_{1\leq j\leq m} \frac{1}{n}\sum_{1}^{n} f(x_i,\theta_j,\rho_{\theta_j}) \leq \sup_{1\leq j\leq m} \mu(\theta_j) + \varepsilon\} = 1,$$

so, $p\{\limsup_{n\to\infty} \sup_{\underline{\theta}\in\Theta} \frac{1}{n}\sum_{1}^{n} W(x_i,\theta) \leq \sup_{\underline{\theta}\in\Theta} \mu(\theta) + \varepsilon\} = 1$

since it is true for all $\varepsilon > 0$, it is also true for $\varepsilon = 0$.

$$\therefore p\{\limsup_{n\to\infty} \sup_{\underline{\theta}\in\Theta} \frac{1}{n}\sum_{1}^{n} W(x_i,\theta) \leq \sup_{\underline{\theta}\in\Theta} \mu(\theta)\} = 1.$$

Now, let $\rho > 0$, and $S = \{\theta : |\theta - \theta_{\circ}| \geq \rho\}$, then $S$ is compact and

$$p\{\limsup_{\theta_\circ} \sup_{\substack{n\to\infty \\ \theta\in S}} \frac{1}{n}\sum_1^n W(x_i,\theta) \le \sup_{\theta\in S}\mu(\theta)\} = 1.$$

By Shannon-kolmogorov Information Inequality, we have

$$\mu(\theta) = -K(\theta_\circ,\theta) = \int W(x,\theta)f(x,\theta_\circ)dv(x) < 0, \qquad \text{for } \theta \in S,$$

since $W(\mathrm{x},\theta)$ is bounded above by an integrable function $K(x)$

and by using the Fatou-Lebesgue Theorem,

$$\limsup_{\theta'\to\theta}\mu(\theta') = \limsup_{\theta'\to\theta} E\,(W(\mathrm{x},\theta'))$$

$$\le E\limsup_{\theta'\to\theta} W(\mathrm{x},\theta')$$

$$\le E(W(\mathrm{x},\theta)) = \mu(\theta) \quad \text{for every} \quad \underline{\theta}\in\Theta.$$

So $\mu(\theta)$ is upper semi continuous and hence achieves its

maximum value on *S*.

Let $\quad \delta = \sup_{\theta\in S}\mu(\theta), \text{ so } \delta < 0 \quad$ and

$$p\{\limsup_{\theta_\circ} \sup_{\substack{n\to\infty \\ \theta\in S}} \frac{1}{n}\sum_1^n W(x_i,\theta) \le \delta\} = 1.$$

Thus, with probability 1, there exists an N such that for all *n>N*,

$$\sup_{\theta\in S} \frac{1}{n}\sum_1^n W(\mathrm{x}_i,\theta) \le \delta/2 < 0, \text{ say.}$$

But

$$\frac{1}{n}\sum_1^n W(\mathrm{x}_i,\hat\theta_n) = \sup_{\theta\in\Theta}\frac{1}{n}\sum_1^n W(\mathrm{x}_i,\theta) \ge 0,$$

the sum is equal to 0 for $\quad \theta = \theta_\circ$.

This implies that $\quad \hat\theta_n \notin S \quad$ for $\quad n > N$, that is, $|\hat\theta_n - \theta| < \rho$.

Since $\rho$ is arbitrary, we get

$$\lim_{n\to\infty} p\{|\hat\theta_n - \theta| < \varepsilon\} = 1.$$

So, the sequence of MLE $\hat\theta_n$ of $\theta$ is strongly consistent.

i.e $\qquad \hat\theta_n \xrightarrow{\text{a.s}} \theta_\circ$.

## 2.2- Asymptotic Normality of MLE

The basis of large-sample tests and confidence intervals is the general property that the MLE has a limiting normal distribution around the true parameter value as mean and with a variance that easily calculated. Indeed, this is a much stronger result than consistency and requires additional regularity conditions.

## Theorem (2.2.1)

Let $x_1$, $x_2$,… be i.i.d with density $f(x/\underline{\theta})$ (with respect to $dv$), and let $\underline{\theta}_\circ$ denote the true value of the parameter. We make the following regularity conditions.

1. $\Theta$ is an open subset of $R^k$,

2. second partial derivatives of $f(x/\underline{\theta})$ with respect to $\underline{\theta}$ exist and are continuous for all $x$, and may be passed under the integral sign in $\int f(x/\underline{\theta})dv(x)$,

3. there exists a function $K(x)$ such that $E_{\underline{\theta}_\circ}(K(x)) < \infty$ and each component of $\psi'(x,\underline{\theta})$ is bounded in absolute value by $K(x)$ uniformly in some neighborhood of $\theta_\circ$,

4. $B(\underline{\theta}_\circ) = -E_{\underline{\theta}_\circ}[\psi'(x,\underline{\theta}_\circ)]$ is positive definite, the information matrix for a single observation,

5. $f(x/\underline{\theta}) = f(x/\underline{\theta}_\circ)$ (a.e $dv$) $\Rightarrow \underline{\theta} = \underline{\theta}_\circ$,

then, there exists a strongly consistent sequence $\hat{\underline{\theta}}_n$ of roots of the likelihood equation such that

$$\sqrt{n}(\hat{\underline{\theta}}_n - \underline{\theta}_\circ) \xrightarrow{L} N(0, B(\underline{\theta}_\circ)^{-1}).$$

## Proof:

Let $S_\rho = \{\theta : |\theta - \theta_\circ| \le \rho\}$, for some $\rho > 0$, be a compact neighborhood of $\theta_\circ$ on which components of $\psi'(x,\theta)$ are uniformly bounded by $K(x)$ as in (3). The existence of a strongly consistent sequence $\hat\theta_n$ of solutions of $\dfrac{\partial \ell_n(\theta)}{\partial\theta} = 0$ follows from theorem (2.1.1) with $\Theta = S_\rho$. Conditions (1), (2) and (5) of that theorem (2.1.1) are already satisfied. Condition (4) follows from continuity of $f(x/\theta)$ in $\theta$.

For checking condition (3), we expand $W(x,\theta)$ about $\theta_\circ$ using Taylor expansion as

$$W(x,\theta) = w(x,\theta_\circ) + \psi(x,\theta_\circ)^T(\theta - \theta_\circ) + (\theta - \theta_\circ)^T \int_0^1\int_0^1 \lambda\psi'(x,\theta_\circ + u\lambda(\theta - \theta_\circ))\,du\,d\lambda\,(\theta - \theta_\circ).$$

Since $W(x,\theta_\circ) = 0, \psi(x,\theta_\circ)$ is integrable, and the components of $\psi'$ are bounded by $K(x)$ uniformly on $S_\rho$, hence $W(x,\theta)$ is bounded uniformly on $S_\rho$ by an integrable function. We have

$$D_\theta \ell_n(\theta) = \sum_1^n \psi(x_i,\theta), \text{ where } D_\theta \text{ is the differential operator } \frac{\partial}{\partial\theta}.$$

By Taylor theorem, we expand $D_\theta \ell_n$ about $\theta_\circ$ we get

$$D_\theta \ell_n(\theta) = D_\theta \ell_n(\theta_\circ) + \int_0^1 \sum_1^n \psi'(x_i,\theta_\circ + \lambda(\theta - \theta_\circ))\,d\lambda(\theta - \theta_\circ)$$

$$D_\theta \ell_n(\theta) = D_\theta \ell_n(\theta_\circ) + \int_0^1 \sum_1^n \psi'(x_i,\theta_\circ + \lambda(\theta - \theta_\circ))\,d\lambda(\theta - \theta_\circ).$$

Now let $\theta = \hat\theta_n$ where $\hat\theta_n$ is any strongly consistent sequence satisfying

$$D_\theta \ell_n(\hat\theta_n) = 0 \quad \text{and divide by } \sqrt{n}$$

$$\frac{1}{\sqrt{n}} D_\theta \ell_n(\theta_\circ) = A_n \sqrt{n}(\hat\theta_n - \theta_\circ),$$

where

$$A_n = -\int_0^1 \frac{1}{n} \sum_1^n \psi'(x_i,\theta_\circ + \lambda(\hat\theta_n - \theta_\circ))\,d\lambda.$$

Because $E_{\theta_\circ}(\psi(x,\theta_\circ)) = 0$ and $Var_{\theta_\circ}(\psi(x,\theta_\circ)) = B(\theta_\circ)$, hence by the central limit theorem, we obtain that

19

$$\frac{1}{\sqrt{n}} D_\theta \ell_n(\theta_\circ) = \sqrt{n}(\frac{1}{n}\sum_1^n \psi(x_i,\theta_\circ)) \overset{L}{\to} Z \quad N(0, B(\theta_\circ)).$$

Now, let $\varepsilon > 0$. Since $E_{\theta_\circ}(\psi'(x,\theta))$ is continuous in $\theta$ from condition (3), so there is a $\rho > 0$ such that

$$|\theta - \theta_\circ| < \rho \text{ implies } |E_{\theta_\circ}(\psi'(x,\theta)) + B(\theta_\circ)| < \varepsilon.$$

From the SLLN, that with probability 1 there is an integer N such that $n > N$ implies

$$\sup_{\theta \in S_\rho} |\frac{1}{n}\sum_1^n \psi'(x_i,\theta) - E_{\theta_\circ}(\psi'(x,\theta))| < \varepsilon.$$

Then assuming N is so large that $n > N$ implies

$$|\hat\theta_n - \theta_\circ| < \rho, \quad n > N \text{ implies}$$

$$|A_n - B(\theta_\circ)| \le \int_0^1 |\frac{1}{n}\sum_1^n \psi'(x_i, \theta_\circ + \lambda(\hat\theta_n - \theta_\circ)) + B(\theta_\circ)| d\lambda$$

$$\le \int_0^1 \sup_{\theta \in S_\rho} [|\frac{1}{n}\sum_1^n \psi'(x_i,\theta) - E_{\theta_\circ}(\psi'(x,\theta))| +$$

$$|E_{\theta_\circ}(\psi'(x,\theta)) + B(\theta_\circ)|] d\lambda \le 2\varepsilon.$$

Hence $A_n \overset{a.s}{\to} B(\theta_\circ)$ implies $A_n^{-1}$ will exist and by Slutsky's theorem

$$\sqrt{n}(\hat\theta_n - \theta_\circ) = A_n^{-1}\frac{1}{\sqrt{n}}D_\theta \ell_n(\theta_\circ) \overset{L}{\to} B(\theta_\circ)^{-1}Z \quad N(0, B(\theta_\circ)^{-1}). \text{ Hence,}$$

under the conditions of theorem (3.2.1) if there is a unique solution of the likelihood equation for every $n$, this sequence of solutions will be consistent and asymptotically normal.

## 2.3- Asymptotic Efficiency of MLE

Under the conditions of theorem (2.2.1), the MLE $\hat\theta_n$ was seen to be asymptotically unbiased in a reasonably strong sense, because whatever be the true value of $\theta$, $\sqrt{n}(\hat\theta_n - \theta)$ is asymptotically normal with mean zero. Moreover, the asymptotic variance of the MLE is $\frac{1}{n}B(\theta)^{-1}$ which is the

Crame'r-Rao lower bound for the variance of any unbiased estimate of $\underline{\theta}$ based on a sample of size $n$.

## *Definition (2.3.1)*

Let $x_1$, $x_2$,…be i.i.d random variables with distribution depending upon a parameter $\underline{\theta} \in \Theta$.

A sequence of estimators $\{\underline{\tilde{\theta}}_n\}$ of $\underline{\theta}$, with $\underline{\tilde{\theta}}_n$ is a function of $x_1$,…, $x_n$, such that $\quad \sqrt{n}(\underline{\tilde{\theta}}_n - \underline{\theta}) \overset{L}{\to} N(0, \text{var}(\underline{\tilde{\theta}}))$.

Whatever be the true value of $\underline{\theta}$, is said to be asymptotically efficient if $\text{var}(\underline{\tilde{\theta}}) = B(\theta)^{-1}$ for all $\underline{\theta} \in \Theta$.

Note that by definition, the MLE is asymptotically efficient under the conditions of theorem (3.2.1).

## 2.4-Asymptotic Sufficiency of MLE

An explained about the asymptotic theory, one way of explaining the good asymptotic properties of MLE's is via their asymptotic sufficiency. The proof is implicit in the earlier arguments, under the usual regularity conditions of that theorems.

The multiparameter case is immediate from the relationship:

$$\sqrt{n}(\underline{\hat{\theta}} - \underline{\theta})B(\theta) = \frac{1}{\sqrt{n}} \sum_{1}^{n} \psi(x_i, \underline{\theta})......(*)$$

such that, the relationship (*) between the MLE and the efficient score vector, the asymptotic sufficiency can also be expressed in terms of $\psi(x, \underline{\theta})$.

## 2.5-Restricted Maximum Likelihood Estimators

In certain problems under one of possible number of assumptions, the vector parameter $\underline{\theta} \in \Theta$ is further restricted to some subspace H of $\Theta$. In that event if we wish to estimate $\underline{\theta}$ under such assumptions we shall need to maximize $f(x/\underline{\theta})$ within H instead of $\Theta$.

If $\hat{\underline{\theta}}$ is that value that maximizes $f(x/\underline{\theta})$ within H

i.e $f(x,\hat{\underline{\theta}}) = \max\limits_{\underline{\theta} \in H} f(x/\underline{\theta})$.

Then $\hat{\underline{\theta}}$ is called the restricted MLE (RMLE) of $\underline{\theta}$.

Usually, the restricted space H is described by a set of side conditions, i.e $\underline{\theta} \in H$ if $\underline{\theta}$ satisfies the side equations:

$$h_1(\underline{\theta}) = 0, h_2(\underline{\theta}) = 0, \ldots, h_r(\underline{\theta}) = 0, \text{i.e } \underline{h}(\underline{\theta}) = \underline{0}$$

One way of finding the RMLE is to use the method of Lagrange multipliers. This method tells us that $\hat{\underline{\theta}}$ satisfies:

$$\left. \begin{aligned} \frac{\partial}{\partial \theta_i} \ell_n(\underline{\theta}) + \sum_{j=1}^{r} \lambda_j \frac{\partial}{\partial \theta_i} h_j(\underline{\theta}) = 0 \\ h_j(\underline{\theta}) = 0 \end{aligned} \right\} \quad \begin{aligned} i &= 1,2,..,k \\ j &= 1,2,..,r \end{aligned} \quad (r < k)$$

In matrix notation:

$$D_{\underline{\theta}} \ell_n(\underline{\theta}) + \mathrm{H}(\underline{\theta})\underline{\lambda} = \underline{0}$$
$$\underline{h}(\underline{\theta}) = \underline{0}$$

where $D_{\underline{\theta}} = \dfrac{\partial}{\partial \underline{\theta}}$, and

$$H(\underline{\theta}) = D_{\underline{\theta}} \, h^{T}(\underline{\theta}) = \begin{pmatrix} \dfrac{\partial}{\partial \theta_1} \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \dfrac{\partial}{\partial \theta_k} \end{pmatrix} (h_1(\underline{\theta}),....,h_r(\underline{\theta}))$$

i.e

$[H(\underline{\theta})]_{i,j} = \dfrac{\partial}{\partial \theta_i} h_j(\underline{\theta})$ and $\underline{\lambda} = \begin{pmatrix} \lambda_1 \\ \cdot \\ \cdot \\ \cdot \\ \lambda_r \end{pmatrix}$ are the Lagrange multipliers.

We assume that rank $(H(\underline{\theta})) = r$.

## *Theorem (2.5.1)*

Let $\underline{\theta}$ be an identifiable *k*-dimensional parameter of a distribution and let $\hat{\underline{\theta}}(x)$ be the RMLE of $\underline{\theta}$ based on a random sample x from the distribution under the side conditions $h(\underline{\theta}) = 0, \underline{h}$ being *r*-dimensional. Then under the regularity conditions of theorem (3.1.1).

(i) $\hat{\underline{\theta}}$ is strongly consistent (i.e $\hat{\underline{\theta}} \xrightarrow{a.s} \underline{\theta}$)

(ii) $\hat{\underline{\theta}}(\underline{x}) \sim N(\underline{\theta}, \dfrac{p(\underline{\theta})}{n})$ for large sample size *n*, where matrix $p(\underline{\theta})$ satisfies:

$$\begin{pmatrix} B(\underline{\theta}) & H(\underline{\theta}) \\ H^{T}(\underline{\theta}) & 0 \end{pmatrix}^{-1} = \begin{pmatrix} p(\underline{\theta}) & Q(\underline{\theta}) \\ Q^{T}(\underline{\theta}) & R(\underline{\theta}) \end{pmatrix}$$

i.e $\quad p(\underline{\theta})B(\underline{\theta})p(\underline{\theta}) = p(\underline{\theta})$.

## *Proof:*

(i) For any $\underline{\theta} \in H$, (i.e $h(\underline{\theta}) = 0$) we have

$$Z(\underline{\theta}') < Z(\underline{\theta}) \text{ where } Z = E(\log L_n(\underline{\theta})) = E(\ell_n(\underline{\theta})).$$

If $\underline{\theta}' \neq \underline{\theta}$, and $\dfrac{1}{n}\ell_n(\underline{\theta}') \xrightarrow{a.s} Z(\underline{\theta}')$, so that $\underline{\hat{\theta}}$, the maximizing

value of $\dfrac{1}{n}\ell_n(\underline{\theta}')$ is close to $\underline{\theta}$, the maximizing point of $Z(\underline{\theta}')$.

Since $\underline{\theta} \in H$, it follows that $\hat{\underline{\theta}}$ is close to H, in which case

the maximizing value, $\hat{\underline{\hat{\theta}}}$ of $\dfrac{1}{n}\ell_n(\underline{\theta}')$ in H, must close to $\hat{\underline{\theta}}$

and hence close to $\underline{\theta}$.

Under the same regularity conditions of theorem (3.1.1)
$\hat{\underline{\hat{\theta}}}(x) \xrightarrow{a.s} \underline{\theta}$ as $n \to \infty$ establishing part (i).

(ii) since $\hat{\underline{\hat{\theta}}}$ satisfies

$$\left. \begin{array}{r} D_\theta \ell_n(\hat{\underline{\hat{\theta}}}) + \mathrm{H}(\hat{\underline{\hat{\theta}}})\,\lambda = 0 \\[2mm] h(\hat{\underline{\hat{\theta}}}) = 0 \end{array} \right\} \qquad \ldots\ (1)$$

We have by Taylor's expansion and the fact that $(\hat{\underline{\hat{\theta}}})$ is

consistent:

$$\left. \begin{array}{r} D_\theta \ell_n(\underline{\theta}) + D_\theta^2 \ell_n(\underline{\theta})\,(\hat{\underline{\hat{\theta}}} - \underline{\theta}) + \mathrm{H}(\underline{\theta})\,\lambda + D_\theta(\mathrm{H}(\underline{\theta})\,\lambda)(\hat{\underline{\hat{\theta}}} - \underline{\theta}) \approx 0 \\[3mm] h(\underline{\theta}) + D_\theta\, h(\underline{\theta})(\hat{\underline{\hat{\theta}}} - \underline{\theta}) \approx 0 \\[3mm] \text{i.e}\quad \mathrm{H}^T(\underline{\theta})\,(\hat{\underline{\hat{\theta}}} - \underline{\theta}) \approx 0 \\[3mm] \ldots\ldots\ldots(2) \end{array} \right\}$$

($\approx$ means approximately)

Also from (i), we have [expand the first term around $(\hat{\underline{\theta}})$]

$$D_\theta \ell_n(\hat{\underline{\theta}}) + D_\theta^2 \ell_n(\hat{\underline{\theta}})\,(\hat{\underline{\hat{\theta}}} - \hat{\underline{\theta}}) + \mathrm{H}(\hat{\underline{\hat{\theta}}})\,\lambda \approx 0.$$

That is,

$$D_\theta^2 \ell_n(\hat{\underline{\theta}})\,(\hat{\underline{\hat{\theta}}} - \hat{\underline{\theta}}) + \mathrm{H}(\hat{\underline{\hat{\theta}}})\,\lambda \approx 0,$$

thus

$| \mathrm{H}(\hat{\hat{\theta}}) \, \underline{\lambda} \,|=| \hat{\hat{\theta}} - \hat{\theta} \,|= 0 \,| \hat{\hat{\theta}} - \underline{\theta} \,|$. So that

$| \mathrm{H}(\underline{\theta}) \, \underline{\lambda} \,| = 0 \,| \hat{\hat{\theta}} - \underline{\theta} \,|$.

And  $| \mathrm{H}(\underline{\theta}) \, \underline{\lambda}(\hat{\hat{\theta}} - \underline{\theta}) \,|= 0 \,| \hat{\hat{\theta}} - \underline{\theta} \,|^2$  and therefore can be omitted from (2) to obtain

$$D_{\underline{\theta}}\ell_n(\underline{\theta}) + D^2_{\underline{\theta}}\ell_n(\underline{\theta})\,(\hat{\hat{\theta}} - \underline{\theta}) + \mathrm{H}(\underline{\theta})\,\underline{\lambda} \approx 0$$

$$\mathrm{H}^T(\underline{\theta})\,(\hat{\hat{\theta}} - \underline{\theta}) \approx 0$$

or

$$\begin{pmatrix} \dfrac{-1}{n}D^2_{\underline{\theta}}\ell_n(\underline{\theta}) & \mathrm{H}(\underline{\theta}) \\ \mathrm{H}^T(\underline{\theta}) & 0 \end{pmatrix} \begin{pmatrix} \sqrt{n}(\hat{\hat{\theta}} - \underline{\theta}) \\ -\dfrac{1}{\sqrt{n}}\underline{\lambda} \end{pmatrix} = \begin{pmatrix} \dfrac{1}{\sqrt{n}}D_{\underline{\theta}}\ell_n(\underline{\theta}) \\ 0 \end{pmatrix}.$$

But

$$\frac{1}{\sqrt{n}}D_{\underline{\theta}}\ell_n(\underline{\theta}) \xrightarrow{L} N(\underline{0}, B(\underline{\theta})), \text{ and}$$

$$-\frac{1}{n}D^2_{\underline{\theta}}\ell_n(\underline{\theta}) \xrightarrow{a.s} B(\underline{\theta})$$

Hence for large $n$

$$\begin{pmatrix} \sqrt{n}(\hat{\hat{\theta}} - \underline{\theta}) \\ -\dfrac{1}{\sqrt{n}}\underline{\lambda} \end{pmatrix} \approx \begin{pmatrix} B(\underline{\theta}) & \mathrm{H}(\underline{\theta}) \\ \mathrm{H}^T(\underline{\theta}) & 0 \end{pmatrix}^{-1} \begin{pmatrix} N(\underline{0}, B(\underline{\theta})) \\ \underline{0} \end{pmatrix}$$

or $\begin{pmatrix} \sqrt{n}(\hat{\hat{\theta}} - \underline{\theta}) \\ \dfrac{-1}{\sqrt{n}}\underline{\lambda} \end{pmatrix} \approx \begin{pmatrix} P(\underline{\theta}) & Q(\underline{\theta}) \\ Q^T(\underline{\theta}) & R(\underline{\theta}) \end{pmatrix} \begin{pmatrix} N(\underline{0}, B(\underline{\theta})) \\ \underline{0} \end{pmatrix},$

so,

$$\sqrt{n}(\hat{\hat{\theta}} - \underline{\theta}) \approx P(\underline{\theta})N(\underline{0}, B(\underline{\theta})),$$

then

$$\hat{\hat{\theta}} \sim N(\underline{\theta}, \frac{1}{n}P(\underline{\theta})B(\underline{\theta})P(\underline{\theta})).$$

Since $B(\underline{\theta})$ is symmetric, so is $P(\underline{\theta})$. Also

$$\begin{pmatrix} B(\bar{\underline{\theta}}) & \mathrm{H}(\underline{\theta}) \\ \mathrm{H}^T(\underline{\theta}) & 0 \end{pmatrix} \begin{pmatrix} P(\underline{\theta}) & Q(\underline{\theta}) \\ Q^T(\underline{\theta}) & R(\underline{\theta}) \end{pmatrix} = \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix}. \qquad \dots\dots\dots (3)$$

That is

$$B(\theta)P(\theta) + H(\theta)Q^T(\theta) = I \qquad\qquad \text{.........(4)}$$

$$H^T(\theta)P(\theta) = 0.$$

The latter implies

$$0 = p^T(\theta)H(\theta) = P(\theta)H(\theta).$$

Now, multiply both sides of (4) by $P(\theta)$, we get

$$P(\theta)B(\theta)P(\theta) + P(\theta)H(\theta)Q^T(\theta) = P(\theta)$$

implies

$$P(\theta)B(\theta)P(\theta) = P(\theta)$$

and

$$\hat{\theta} \sim N(\theta, \frac{P(\theta)}{n}) \quad \text{for large } n, \text{ establishing the part (ii).}$$

## 2.6- Confidence Interval for Large Samples

We know that under very general conditions, the MLE $\hat{\theta}$

is distributed normally about mean $\theta$ with variance

$$\text{var}(\hat{\theta}) = \frac{1}{nE[\dfrac{\partial^2}{\partial\theta^2}\log f(x/\theta)]}.$$

For large samples:

$$p\{\hat{\theta} - Z_{\alpha/2}\sqrt{\text{var}(\hat{\theta})} \le \theta \le \hat{\theta} + Z_{\alpha/2}\sqrt{\text{var}(\hat{\theta})}\} = 1 - \alpha$$

where $\alpha$ is the significance level.

# REFERENCES

[1] Ash, R. (1972): "Real Analysis and Probability", Academic Press, New York.

[2] Cox, D.R. and Hinkley, D.V. (1974): "Theoretical Statistics", Chapman and Hall, London.

[3] Dedwicz, E.J. (1976): "Introduction to Statistics and Probability", Holt, Rinehart and Winston, U.S.A.

[4] Ferguson, T.S. (1996): "A Course in Large Sample Theory", Chapman and Hall, Great British. New York.

[5] Klimov, G. (1986): "Probability Theory and Mathematical Statistics" Mir, Publishers. Moscow.

[6] Rao, C.R. (1973): "Linear Statistical Inference and its Applications". John Wiley and Sons, New York.

[7] Sharma, J.N. (1977): "Topology", Meerut College, Meerut.

[8] Silvey, S.D. (1978): "Statistical Inference". Chapman and Hall, London.

[9] Wald A. (1949): "Note on The consistency of the Maximum likelihood Estimate". Ann.Math.Statist. Vol.20, pp 595-601.

[10] Wolfowitz, J. (1949): "On Walds Proof of the Consistency of the Maximum likelihood Estimate" Ann. Math. Statist.Vol.20, pp. 601-602.