## MAXIMUM LIKELIHOOD AND LEAST SQUARES ESTIMATION

by

## Jennifer McDonald DaCosta Goddard

A thesis submitted to the Faculty of Graduate Studies and Research in partial fulfillment of the requirements for the degree of Master of Science

Department of Mathematics McGill University Montreal

.

July, 1978

## MAXIMUM LIKELIHOOD AND LEAST SQUARES ESTIMATION

by

## JENNIFER McDONALD DaCOSTA GODDARD

# RÉSUMÉ

Cette thèse examine deux méthodes d'estimation parmi les plus importantes et les plus souvent utilisées - l'estimation du maximum de vraisemblance pour laquelle on supposée connue la distribution des variables (ou vectors) aléatoires et d'autre part l'estimation des moindres carrés où cette fois la distribution de variables aléatoires n'est pas nécessairement connue.

Quelques définitions et résultats préliminaires seront présentés dans le premier chapitre. Ceux-ci seront utiles dans les chapitres suivants. Le second chapitre portera sur les propriétés souhaitables des estimateurs qui seront rencontrés aux chapitres trois et quatre.

L'estimation du maximum de vraisemblance sera discutée au chapitre trois et l'emphase sera placée sur les propriétés assymptotiques des estimateurs. Le quatrième chapitre traitera de l'estimation des moindres carrés. Dans une section de ce chapitre, le problème d'estimabilité sera étudié. On trouvera de nombreuses références dans la bibliographie.

Department of Mathematics

M.Sc.

### ACKNOWLEDGEMENTS

My most sincere appreciation and thanks to Professor V. Seshadri for his undeniable guidance and support during the preparation of this thesis, and to the Canadian Commonwealth Scholarship Committee for financing my graduate studies.

I thank my fiancée Anthea whose frequent letters of encouragement have helped me greatly. I also thank my brothers Rupert and Horace and their families for providing me with an atmosphere conducive to studying.

Finally, I thank Mrs. Pamela Neblett for her patience, diligence and accuracy in the typing of this thesis.

# TABLE OF CONTENTS

.

CHAPTER 1	Preliminary Results and Definitions 1
CHAPTER 2	Desirable Properties of Estimators 10
2.1	Introduction 10
2.2	The Problem 10
2.3	Consistency 12
2.4	Sufficiency 13
2.5	Unbiasedness 19
2.6	Minimum Variance 20
2.7	Completeness 23
2.8	Efficiency 31
CHAPTER 3	Maximum Likelihood Estimation
3.1	Introduction
3.2	Regularity Assumptions
3.3	Properties of the Score Vector
3.4	Properties of the Information Matrix 39
3.5	Iterative methods
3.6	Properties of MLE's 47
CHAPTER 4	Least Squares Estimation
4.1	Introduction
4.2	Asymptotic Properties of Least Squares 59
4.3	Robustness of the method of Least Squares.65

4.4	A Different Computational Technique	66
4.5	Estimability of Linear Contrasts	73
	Bibliography	83

•

Page

- 1.15 Theorem: If A is a symmetric matrix, then there exists an orthogonal matrix P such that PAP' =  $\Lambda$ , where  $\Lambda$  is the diagonal matrix of characteristic roots of A and P the matrix of characteristic vectors.
- 1.16 The singular values of the matrix B (m x n) denoted by Sg(B) are defined to be the positive square roots of the characteristic roots of the matrix B'B, i.e., Sg(B) = +  $\left[ CH(B'B) \right]^{\frac{1}{2}}$ .
- 1.17 Let B (m x n) be an arbitrary real matrix. The rank of B, denoted by r (B) is the maximal number of linearly independent rows (columns) of B.
- 1.18 If  $n \le m$  and r(B) = n, then  $B(m \ge n)$  is said to have full column rank, if  $m \le n$  and r(B) = m, B is said to have full row rank.
- 1.19 The matrix X (n x m) is said to be a right inverse of A (m x n) if A X =  $I_m$ . If Y (n xm) is such that YA =  $I_n$ , then Y is a left inverse of A.
- 1.20 Lemma: If A (m x n) is such that r(A) = n, then there exists a matrix X (n x m), called a right inverse of A, such that A X = I<sub>m</sub>; one choice of X is  $A'(AA')^{-1}$ . If r(A) = m, then there exists a matrix Y such that  $YA = I_n$ ; one choice of Y is  $(A'A)^{-1}A'$ .

1.21 A generalized inverse of A  $(m \times n)$  is a matrix denoted by A<sup>-</sup> of order  $n \times m$  such that

5

 $A \overline{A} = A$ .

- 1.22 The Moore-Penrose generalized inverse of the matrix A (m x n) is a matrix G (n x m) such that the following four conditions are satisfied:
  - (1) A G A = A
  - $(2) \quad G \land G = G$
  - (3) (AG)' = AG
  - (4) (GA)' = GA.

The matrix G is called a  $g_{1234}$  (A) and is denoted by A<sup>+</sup>. If condition (1) is satisfied, G is called a weak generalized inverse of A and is denoted by  $g_1(A)$ . If conditions (1) and (2) are satisfied G is called a symmetric generalized inverse of A and is denoted by  $g_{12}(A)$ . If conditions (1) and (3) are satisfied G is called a least squares generalized inverse and is denoted by

$$g_{13}(A)$$
 or  $A_{\ell s}^{-}$ .

1.23 Theorem: Let A (m x n) be an arbitrary real matrix; then there exists orthogonal matrices P (m x m) and Q (n x n) such that  $\begin{bmatrix} n & n \end{bmatrix}$ 

$$P'AQ = \begin{bmatrix} Dr & 0\\ 0 & 0 \end{bmatrix}$$
, where  $D_r$  is the

diagonal matrix of singular values of A, P the matrix of

1.24 Lemma:  $A^+$  and  $A^-_{\&s}$  are unique.

1.25 Lemma: Using (1.23) we note that one representation of  $A^+$  is

$$A^{+} = Q \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} P'$$

1.26 Theorem: For comformable matrices A and B and for any choice of the generalized inverses A and B, we have

$$r(A,B) = r(A) + r([I-AA^{-}]B),$$
 (1.26.1)

$$= r(B) + r([I-BB]A). \qquad (1.26.2)$$

$$r\begin{bmatrix} A\\ B\end{bmatrix} = r(A) + r([B(I-A^{-}A)]), \qquad (1.26.3)$$

$$= r(B) + r([A(I-B^{B})]). \qquad (1.26.4)$$

1.27 Lemma. 
$$r(AB) \leq r(A)$$
, (1.27.1)

$$r(AB) \leq r(B)$$
. (1.27.2)

1.28 Lemma: Let the matrix C (m x n),  $n \leq m$ , have full column rank n, and the matrix R (n x m),  $n \leq m$ , have full row rank n. Then for any matrix A (n x m), we have

$$r(A) = r(CA) = r(AR).$$

1.29 The trace of a matrix A (n x n) denoted by tr(A) is defined to be the sum of the diagonal elements; alternatively the trace of A is the sum of the characteristic roots of A.

7

1.30 If A (n x n) is an idempotent matrix of rank r, Lemma: then tr(A) = r,

Lemma: Let  $\alpha > 0$  and A (m x n) be a real matrix. 1.31 Then  $\lim_{\alpha \to 0} \left[ A'A + \alpha I \right]^{-1} A' = (A'A)^{-1} A'.$ 

Proof: by (1.23), we have

A	=	$P\begin{bmatrix} D\\ 0 \end{bmatrix}$	0]0'	
---	---	----------------------------------------	------	--

hence

and

 $A'A = Q \begin{bmatrix} D^2 & 0 \\ r & 0 \end{bmatrix} Q'$ (1.31.1) $A'A + \alpha I_{n} = Q \begin{bmatrix} D_{r} + \alpha I_{r} & 0 \\ 0 & \alpha I_{n-r} \end{bmatrix} Q'.$ (1.31.2)

We note that for all  $\alpha > 0$ , A'A +  $\alpha I_n$  is invertible.

Thus

where

 $\begin{bmatrix} A'A + \alpha I_n \end{bmatrix}^{-1} = Q \begin{bmatrix} \Lambda(\alpha) & 0 \\ 0 & \alpha^{-1}I \end{bmatrix} Q'$ (1.31.3)  $\Lambda(\alpha) = \left[ D_r + \alpha I_r \right]^{-1},$ hence,  $\begin{bmatrix} A'A + \alpha I_n \end{bmatrix}^{-1}A' = Q \begin{bmatrix} \Lambda(\alpha) & 0 \\ 0 & \alpha^{-1}I_{n-r} \end{bmatrix} Q'Q \begin{bmatrix} D_r & 0 \\ 0 & 0 \end{bmatrix} P'$  $= Q \begin{bmatrix} \Lambda(\alpha)D_{r} & 0 \\ 0 & 0 \end{bmatrix} P', \quad (1.31.4)$ 

but  $\lim_{\alpha \to 0} \Lambda(\alpha) = D_r^{-2}$ ,

hence

$$\lim_{\alpha \to 0} \left[ A'A + \alpha I_n \right]^{-1} A' = Q \begin{bmatrix} D_r^{-1} & 0 \\ 0 & 0 \end{bmatrix} P'$$
$$= (A'A)^{-} A' \cdot$$

Note: The above lemma can be generalized by replacing  $\alpha I_n$  by D( $\alpha$ ) a diagonal matrix of positive, not necessarily identical elements.

8

1.32

The density function f(x) is a member of the exponential family if

 $f(x) = A(x) B(\theta) \exp [C(\theta) D(x)]$ , a < x < b, where

(1)  $\theta$  is the vector of parameters,

- (2) a or b does not depend on  $\theta_{\mathbf{y}}$
- (3) A(x) and D(x) are continuous for a < x < b,
- (4) C ( $\theta$ ) is a non-trivial continuous function of  $\theta$ .

Examples of Exponential families are

(1) { f(x): f(x) = 
$$\frac{1}{\theta} \exp(-\theta x), x \ge 0, \theta > 0$$
}  
(11) { f(x): f(x) =  $\frac{1}{\theta} e^n x^{n-1} \exp(-\theta x), x \ge 0, \theta > 0, n > 0$ }

1.33 Two subspaces are said to be virtually disjoint if the null vector is the only vector they have in common.

1.34 Let  $A = \begin{bmatrix} E & F \\ G & H \end{bmatrix}$  with H nonsingular.

The Schur complement of E in A denoted by A/E is defined as;

$$A/E = E - F H^{-1} G,$$

#### CHAPTER II

## DESIRABLE PROPERTIES OF ESTIMATORS

- 2.1 <u>Introduction</u>. The theory of mathematical estimation was founded by R. A. Fisher in a series of Papers (1921, 1925, 1934 and others). We shall discuss some of the important ideas developed by Fisher in this and a subsequent Chapter. We shall also classify and study the properties of various estimators and give some examples to strengthen the discussion.
- 2.2 <u>The Problem</u>. Suppose we have obtained a random sample from a population, whose distribution has a known mathematical form, but involves a certain number of unknown parameters. We would like to find estimates for the values of these unknown parameters in the parent population based on the information given by the sample.

In general, there will be an infinity of functions of the sample values that might be proposed as estimates of the parameters. We are faced with the following question; "how should we best combine the data to form these estimates?" This question leads to another, namely, "what do we mean by "best" estimates?"

We might answer impulsively that the best estimates will be the ones which are nearest to the true values of the parameters to be estimated. However, it must be borne in mind that every suggested estimate is a function of the sample values, and hence must be regarded as an observed value of a certain random variable. We therefore see that we have no means of predicting the individual values assumed by the estimates in any given case, so that the goodness (in the sense of being best) of an estimate cannot be judged from individual values, but only from the distribution of the values which the estimate will assume in the long run, that is, from its sampling distribution.

When most of the "mass" in this distribution is concentrated in some small neighbourhood of the true value, there is a very small probability that the estimate will differ from the true value by more than a small amount. We express this probabilistically as  $P[|t_n - \theta| > \varepsilon] < \delta$ , where  $t_n$  is the estimator of the parameter  $\theta$ , and  $\delta$  is a small positive quantity depending on n and  $\varepsilon$ . From this point of view, an estimate will be "better" in the same measure as its sampling distribution shows a greater concentration about the true value.

We now list some desirable properties, some or all of which we would like our estimator to have:

- (a) Consistency,
- (b) Sufficiency,
- (c) Unbiasedness,
- (d) Minimum variance,
- (e) Completeness,

## (f) Efficiency.

We now define and give examples of statistics which satisfy the properties listed above.

## 2.3. Consistency

<u>Definition</u>: An estimator  $t_n$  computed from a random sample of n values is said to be a consistent estimator for the parameter  $\theta$ , if for any given  $\varepsilon > 0$  there exists an N and a  $\delta$  such that

$$P(|t_{-} \theta| \leq \varepsilon) \geq 1 - \delta$$
 (1.3.1)

for all  $n \ge N$ ; t<sub>n</sub> is then said to converge in probability or to converge stochastically to  $\theta$ .

Note: As defined, consistency is a large sample property.

Example

Consider dF(x) = 
$$\frac{1}{\sqrt{2\pi}} \exp -(\frac{x-\theta}{2})^2$$
,  $-\infty < x < \infty$ 

and let  $x = (x_1, x_2, \dots, x_n)$  be a random sample of size n taken from this population. Then the sample mean  $t = \sum_{i=1}^{n} \frac{x_i}{n}$  is consistent

for the parameter  $\theta$ , since, given  $\epsilon > 0$ , and using the fact that  $\sqrt{n} (t-\theta) \sim N(0,1)$ , we have

$$P(\sqrt{n} | t - \theta| < \varepsilon \sqrt{n}) = \int_{-\varepsilon \sqrt{n}}^{\varepsilon \sqrt{n}} dF(x)$$
$$= \int_{-\varepsilon \sqrt{n}}^{\varepsilon \sqrt{n}} \exp(-\frac{t^2}{2}) \frac{dt}{\sqrt{2\pi}}$$

Clearly this probability can be made as close to one (1) as we please, by making n sufficiently large. For example, if  $\varepsilon = 0.1$  and n = 400, we have,

$$P(20|t-\theta| \leq 2) = P(|t-\theta| \leq .1)$$

# = 0.9544 .

It can be shown that the sample median,  $x_{\frac{1}{2}}$ , is also consistent for the parameter  $\theta$  in the normal population and thus we note that consistent estimators are not necessarily unique. In cases where there are several consistent estimators it is usual to choose the estimator with the smallest variance.

In the case of a random sample of size n taken from a normal population with mean/median  $\mu$  and variance  $\sigma^2$ , it can be shown that the sample mean/median are both asymptotically normal with common expected value  $\mu$  and variances  $\frac{\sigma}{n}$  and  $\frac{\pi\sigma}{2n}$  respectively. Thus we see that the sample mean will be preferred to the sample median as an estimator for the population mean, since the sample mean has asymptotically smaller variance than the sample median.

## 2.4. Suffciency

Many decision theoretic problems can be significantly simplified by a suitable reduction in the complexity of the data. A suitable form of reduction called the criterion of sufficiency was introduced by R. A. Fisher (1925). There he defined a sufficient statistic as one which effectively summarizes all the relevant information supplied by the data.

The following definition offers a small extension

to Fisher's concept of sufficiency.

<u>Definition</u>: A statistic is said to satisfy the criterion of sufficiency when no other statistic which can be calculated from the data provides any additional information as to the value of the parameter to be estimated.

A more rigorous definition is as follows: <u>Definition</u>: Let X denote a random variable (vector) whose distribution depends on a parameter  $\theta \in \Theta$ ,  $\Theta$  being the parameter space. A real (vector) function T of X is said to be sufficient for the parameter  $\theta$  if the conditional distribution of X given T = t is independent of  $\theta$  almost everywhere (t).

The following theorem, due to Fisher, gives us a necessary and sufficient condition for a statistic to be suf-ficient.

### Fisher's Factorization Theorem:

Let X be a random variable whose p.d.f. (p.m.f)  $f(x,\theta)$ depends on the parameter  $\theta \in \Theta$ . A function T = t(x) is sufficient for  $\theta$  if and only if the frequency function factors into a product of the two functions g and h, where g is a function of t(x) and  $\theta$ , and h is a function of x only, i.e.,

 $f(x,\theta) = g(t(x),\theta) h(x)$ .

<u>Proof</u> (if ) Suppose that T = t(x) is sufficient for  $\theta$  then,

$$f(x,\theta) = P_{\theta}(X = x)$$
  
=  $P_{\theta}(X = x, T = t(x))$   
=  $P_{\theta}(T = t(x)) P(X = x | T = t(x))$  (2.4.1)

provided that P(X = x | T = t (x)) exists and is well defined. Hence for the x for which  $f(x, \theta) = 0$  for all  $\theta$ 

we define

$$h(x) = 0$$

and for the x for which  $f(x,\theta) > 0$  for some  $\theta$ we put

$$h(x) = P(X = x | T = t(x))$$
 (2.4.2)

which is independent of  $\theta$ , so that putting

$$g(t(x), \theta) = P_{\Omega} (T = t(x))$$
 (2.4.3)

gives us the required factorization.

Proof (only if) Suppose

$$f(x,\theta) = P_{\theta}(X = x)$$
  
=  $P_{\theta}(T = t(x)) P(X = x | T = t(x))$  (2.4.4)

holds. Fix  $t_0$  for which  $P_{\theta}(T = t_0) > 0$  for some  $\theta \in \bigoplus$ . Then

$$P_{\theta}(X = x | T = t_{0}) = \frac{P_{\theta}(X = x, T = t_{0})}{P_{\theta}(T = t_{0})} \cdot (2.4.5)$$
  
The numerator of (2.4.5) is zero for all  $\theta$  whenever  
t (x)  $\neq$  t\_{0} and is equal to  $P_{\theta}(X = x)$ , whenever t (x) = t\_{0}

We can write the denominator of (2.4.5) as

$$P_{\theta}(T = t_0) = \Sigma \qquad P_{\theta}(X = x)$$

$$\begin{bmatrix} x:t(x) = t_0 \end{bmatrix}$$

$$= \Sigma \qquad g(t(x), \theta) \cdot h(x) \cdot (2.4.6)$$

$$\begin{bmatrix} x:t(x) = t_0 \end{bmatrix}$$

Hence

$$P_{\theta}(\mathbf{X} = \mathbf{x} | \mathbf{T} = \mathbf{t}_{0}) = 0 \qquad \mathbf{t}(\mathbf{x}) \neq \mathbf{t}_{0}$$
$$= \frac{g(t_{0}, \theta) \cdot h(\mathbf{x})}{g(t_{0}, \theta) \cdot h(\mathbf{x}_{1})}, \quad \mathbf{t}(\mathbf{x}) = \mathbf{t}_{0}$$

where 
$$h(x_1) = \Sigma h(x')$$
.  

$$\begin{bmatrix} x':t(x') = t_0 \end{bmatrix}$$

Cancelling  $g(t_0, \theta)$ , we have that,

 $P_{\theta}(X = x | T = t_0)$  is independent of  $\theta$ , for all  $t_0$  and  $\theta$  for which it is defined.

We shall see later that using the criterion of sufficiency in conjunction with that of completeness, we can construct uniformly minimum variance unbiased estimators (UMVUE) for estimable functions. We give here two examples of how to find a sufficient statistic using the Fisher Factorization Criterion. <u>Example I</u>. Let  $x = (x_1, x_2, \dots, x_n)$  be a random sample from a Poisson Population with parameter  $\theta > 0$ . Then the likelihood function is

$$L(\theta | \mathbf{x}) = \frac{n}{11} \frac{\theta^{\mathbf{x}_{\mathbf{i}}}}{\mathbf{x}_{\mathbf{i}}!} \exp(-\theta)$$
$$= \frac{1}{\mathbf{x}_{\mathbf{i}}! \mathbf{x}_{\mathbf{2}}! \cdots \mathbf{x}_{\mathbf{n}}!} \cdot \frac{\theta^{\sum_{i=1}^{n} \mathbf{x}_{\mathbf{i}}}}{\exp(-\mathbf{n}\theta)}$$
$$= h(\mathbf{x}) g(t(\mathbf{x}), \theta)$$

where  $h(x) = \begin{bmatrix} n \\ \Pi \\ i=1 \end{bmatrix}^{-1}$  and  $t(x) = \sum_{i=1}^{n} x_i$ , i=1

thus t (x) is sufficient for  $\theta$ . It is easily seen that  $\frac{t(x)}{n}$  is a consistent and unbiased estimator for  $\theta$ .

Example II. Let  $x = (x_1, x_2, \dots, x_n)$  be a random sample from a distribution with pdf

$$f(x,\theta) = G(\theta)M(x) \quad \text{if } 0 < x < \infty$$
$$= 0 \qquad \text{elsewhere,}$$

where  $\left[G(\theta)\right]^{-1} = \int_{0}^{\theta} M(x) dx$ .

The likelihood function is:

$$L(\theta | \mathbf{x}) = \prod_{i=1}^{n} G(\theta) M(\mathbf{x}_{i}) \quad \text{if all } \mathbf{x}_{i} \in (0, \theta)$$

We can then write  $L(\theta | x)$  as

$$L(\theta | \mathbf{x}) = q(0, \mathbf{y}_{1}) \prod_{i=1}^{n} M(\mathbf{x}_{i}) q(\mathbf{y}_{n}, \theta) \left[G(\theta)\right]^{n}$$

 $\forall x_1, x_2, \dots, x_n$ , where  $y_1$  is the first order statistic,

 $y_n$  is the  $n^{th}$  order statistic, and

q(a,b) = 1 if a < b= 0 if a > b.

If we now let  $h(x) = q(0,y_1) \underset{i=1}{\overset{n}{\amalg}} M(x_i)$ 

and 
$$g(t(x), \theta) = q(y_n, \theta) \cdot [G(\theta)]^n$$
,

we see that the n<sup>th</sup> order statistic  $y_n$  is sufficient for the parameter  $\theta$ .

<u>Note:</u> When  $G(\theta) = \theta^{-1}$  and M(x) = 1 we have the uniform distribution. Another specific example of this form is the distribution with  $M(x) = nx^{n-1}$  and  $G(\theta) = \theta^{-n}$ .

2.5. Unbiasedness

<u>Definition</u>: An estimator  $t_n$  of a parameter  $\theta$  is said to be unbiased if the expected value of  $t_n$  is equal to  $\theta$ , i.e. if  $E[t_n] = \theta$ ; otherwise the estimator is said to be biased with a certain bias, b ( $\theta$ ).

We mentioned earlier that it is very desirable to choose a consistent statistic as an estimator of the population parameter. In general, such estimators are biased for the parameter. However, in some cases it is easy to change these biased estimators into unbiased estimators by the introduction of a constant scale and or shift factor as the case may be. In other instances when the expected value is a very complicated function of the parameter it is not so easy to determine these scale and shift factors. However, Quenouille (1956) has suggested a rather elegant method for overcoming this difficulty in many situations. The interested reader is referred to "The Theory of Advanced Statistics, Vol. II" by M. G. Kendall and A. Stuart.

We now give an example in which a single scale factor is necessary to change a biased estimator into an unbiased estimator.

Example: Let  $x_1, x_2, \ldots, x_n$  be a sample of size n taken from a univariate normal population with known mean  $\mu$  and unknown variance  $\sigma^2$ . Then

$$S = \sum_{i=1}^{\infty} (x_i - \mu)^2 / n$$

is the M.L.E for  $\sigma^2$ . It is however biased, since  $E[S] = \frac{n-1}{n} \sigma^2$ . We therefore see that  $\frac{nS}{n-1}$  is unbiased for  $\sigma^2$ .

## 2.6 Minimum Variance

We mentioned in the discussion of consistent estimators that when there are two or more consistent estimators for a parameter the desirable estimator is the one with the smallest variance. This leads us to a very important class of estimators, namely, minimum variance estimators. The concept of minimum variance has been used since the days of Gauss and Laplace. However, in relatively recent times, Cramér, Rao and Battacharrya have shown that lower bounds for the variance exist. We shall deal with minimum variance estimators which are also consistent and unbiased.

<u>Definition</u>: A consistent estimator  $t_n$  of a parameter  $\theta$  is a minimum variance unbiased estimator if

(i)  $E[t_n] = \theta$ 

and

(ii)  $\lim_{n \to \infty} \operatorname{Var}(t_n) \leq \lim_{n \to \infty} \operatorname{Var}(t_n')$ 

for all consistent unbiased estimators  $t'_n$ . We now state and prove a theorem on minimum variance bounds and discuss methods of determining such bounds. <u>Theorem</u>: If L( $\theta$ ) and t satisfy the regularity assumptions RA1-RA5 and if  $E_{\theta}[t] = \tau(\theta)$ then Var (t)  $\ge \frac{\partial}{\partial \theta} \tau(\theta) J(\theta)^{-1} \frac{\partial}{\partial \theta'} \tau(\theta)$ . (2.6.1)

The univariate analogue of this theorem is known as the Cramér-Rao theorem and the quantity on the right of inequality (2.6.1) is called the Cramér-Rao (sometimes the Frechet-Cramér-Rao) lower bound.

<u>Proof of the Theorem</u>: Let  $X_1$  and  $X_2$  be arbitrary random vectors with means  $\mu_1$  and  $\mu_2$  covariance matrices  $\Sigma_{11}$  and  $\Sigma_{22}$  respectively and cross covariance matrix  $\Sigma_{12}$ , we assume that the matrix  $\Sigma_{22}$  is positive definite.

Using 1.32, we have the following identity

 $\Sigma_{11} = \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} + \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$ (2.6.2)

 $\Sigma = \Sigma = \Sigma = \Sigma$  is the schur complement of 11 = 12 = 22 = 21,

 $\sum_{\substack{11\\11}} \text{ in } \Sigma = \begin{bmatrix} \Sigma & \Sigma \\ 11 & 12 \\ \Sigma & \Sigma \\ 21 & 22 \end{bmatrix}$ 

Setting  $X_1 = t$  and  $X_2 = S$ , the score vector, we have  $\mu_1 = \tau(\theta)$ and  $\mu_2 = 0$ ,

 $\Sigma_{12} = \Sigma_{t,s} = \frac{\partial}{\partial \theta} \tau (\theta) = \Sigma_{21}$  $\Sigma_{22} = \Sigma_{s,s} = J(\theta), \Sigma_{11} = \Sigma_{t,t}.$ 

Thus

$$\Sigma_{11} = \Sigma_{tt} = \frac{\partial}{\partial \theta} \tau \quad (\theta) \quad J(\theta) = \frac{1}{\partial \theta} \tau \quad (\theta) \quad + \quad Var \left[ t - \tau(\theta) - \Sigma_{12} \sum_{22}^{-1} s \right]$$

Hence,

$$\Sigma_{t,t} \geq \frac{\partial}{\partial \theta} \tau(\theta) J(\theta)^{-1} \frac{\partial}{\partial \theta}, \tau(\theta),$$

since a covariance matrix is non-negative definite. We have equality if and only if

$$z = \tau(\theta) + \Sigma_{12} \Sigma_{22}^{-1} S$$

almost everywhere (t). Q.E.D.

<u>Note</u>: For non-negative definite matrices A and B, A  $\ge$  B means that A-B is non-negative definite.

Example: Let  $X \sim N_p(\mu, \Sigma)$ .

t

Suppose a random sample of size n is taken from this normal population and suppose also that  $\Sigma$  is positive definite. Then,

$$\bar{\mathbf{X}} = \sum_{i=1}^{n} \mathbf{x}_{i}/n$$

is sufficient for  $\mu$ , and has minimum variance.

 $\frac{\Pr \circ f}{L(\mu, \Sigma | X)} = (2\pi)^{-\frac{np}{2} | \Sigma |} \sum_{i=1}^{-n/2} \exp \left[ \sum_{i=1}^{n} (x_i - \mu)^{i} \Sigma^{-1} (x_i - \mu) \right].$ 

Therefore,

 $\log L(\mu, \Sigma | X) = -\frac{nP}{2} \log 2\pi - \frac{n}{2} \log |\Sigma| + \frac{n}{2} (x_i - \mu)' \Sigma^{-1} (x_i - \mu),$   $\frac{\partial \ell}{\partial \mu} = \frac{\partial}{\partial \mu} \sum_{i=1}^{n} (x_i - \mu)' \Sigma^{-1} (x_i - \mu)$   $= \sum_{i=1}^{n} \frac{\partial}{\partial \mu} (x_i - \mu)' \Sigma^{-1} (x_i - \mu)$   $= \sum_{i=1}^{n} \left[ -2\Sigma^{-1} x_i + 2\Sigma^{-1} \mu \right]$   $= -2\Sigma^{-1} \left[ \sum_{i=1}^{n} (x_i - \mu) \right]$  $= -2\Sigma^{-1} \left[ \sum_{i=1}^{n} (x_i - \mu) \right].$  Hence,

$$\frac{\partial l}{\partial u} = 0 \implies \bar{x} = \mu$$

Thus  $\overline{\mathbf{x}}$  is the maximum likelihood estimator for  $\mu$  and the Cramér-Rao bound  $\Sigma/n$  is attained by this estimator, a fact easy to verify.

Note that maximum likelihood estimators attain the Cramér-Rao lower bound. Furthermore in order to compute the Cramér-Rao lower bound, it is not necessary to find the variance of the maximum likelihood estimator explicitly. We just set  $\frac{\partial \ell}{\partial \mu} = 0$ , divide by the constant, and the coefficient of  $(t - \mu)$  will be the Cramér-Rao bound; here, t is the maximum likelihood estimator of  $\mu$ .

## 2.7 <u>Completeness</u>

The concept of completeness of a parametric family was introduced by Lehmann and Scheffé (1950). Completeness confers the uniqueness property on the class of sufficient estimators. This property taken in conjunction with that of sufficiency will be used to determine uniformly minimum variance unbiased estimators for estimable functions of a parameter  $\theta$ . <u>Definition</u>: A parametric family of distributions f (x, $\theta$ ), depending on the value of the parameter  $\theta$ , is said to be complete if, for h (x) any statistic independent of the parameter  $\theta$ , the expected value of h(x) being equal to zero implies that h(x) = 0 almost everywhere (x). If the expected value of h(x) = 0 implies that h(x) = 0 only for bounded h(x), then  $f(x, \theta)$  is said to be boundedly complete.

<u>Definition</u>: A statistic  $t = \phi(x)$  is complete, if the family of distribution functions of t is complete.

We list here with proofs a few standard examples of complete families of distribution functions.

2.7.1 The family of binomial distributions with  
p.m.f 
$$\binom{n}{x} p^{x} (1-p)^{n-x}, \quad x = 0, 1, ..., n$$
  
 $0 , is complete, since, if
$$E_{p}\left[h(x)\right] = \sum_{x=0}^{n} h(x) \binom{n}{x} p^{x} (1-p)^{n-x} = 0, \forall p, 0 \le p \le 1.$$$ 

then this n<sup>th</sup> degree polynomial in p must be the unique zero polynomial. That is, the coefficients of each power of p must be zero.

Now,

$$\begin{array}{c} n \\ \Sigma \\ x=0 \end{array} h(x) \binom{n}{x} p^{x} (1-p)^{n-x} \\ n \\ \Sigma \\ x=0 \\ j=0 \end{array} (-1)^{j} \binom{n}{x} \binom{n-x}{j} h(x) p^{x+j} \\ 0 \end{array}$$

We therefore have

=

=

$$1 = p^{0} \Rightarrow x = 0, j = 0$$

hence, h(0) = 0.

 $p^{1} \rightarrow x = 0, j = 1 \text{ or } x = 1, j = 0$ hence  $\binom{n}{1}h(0) + \binom{n}{0}h(1) = 0$ or h(1) = 0, since h(0) = 0.

Continuing in this manner we see that

$$h(0) = h(1) = \dots = h(n) = 0$$
.

2.7.2. The family of all Poisson distributions with

0.E.D.

p.m.f.  $\frac{\theta}{x!} \exp(-\theta)$ ,  $x = 0, 1 \dots \theta > 0$  is complete.

Proof:

$$f(x,\theta) = \frac{\theta}{x^{1}} \cdot \exp(-\theta), \quad x = 0, 1, 2, \dots, \theta > 0.$$

If h(x) is any function of the data, we have

$$E_{\theta}[h(\mathbf{x})] = \sum_{\mathbf{x}=0}^{\Sigma} h(\mathbf{x}) \frac{\theta^{\mathbf{x}}}{\mathbf{x}!} \exp(-\theta) = 0$$

If we expand exp  $(-\theta)$  as a Taylor series about zero, we get

$$E_{\theta}\left[h(x)\right] = \sum_{x=0}^{\infty} h(x) \cdot f(\theta) \equiv 0$$

hence, h(x) = 0 almost everywhere (x), because of the uniqueness of the Taylor series expansion of  $f(\theta)$ .

The family of normal distributions with mean zero and variance  $\sigma^2 > 0$  is not complete, since any odd function h(x) has expectation zero. We now state and prove two theorems, the first of which gives us a means of improving on an arbitrary unbiased estimator while the second theorem gives us two ways in which we can determine a uniformly minimum variance estimator.

### 2.7.3. Completeness and the Rao-Blackwell Theorem

<u>Theorem</u>: If the statistic t is sufficient for the parameter  $\theta$  and  $E_{\theta}[t'] = g(\theta)$ , then,  $t'' = E_{\theta}[t'|t]$ 

is such that t" is an unbiased estimator of g ( $\theta$ ); t" is therefore  $\theta$ -free and is a function of t only,

say, t'' = h(t) such that

Var  $(t'') \leq Var(t')$ 

With equality holding if and only if t' and t are dependent with probability one.

Proof: t sufficient for  $\theta$  implies that the distribution of t' | t is  $\theta$ -free, hence,  $E_{\theta}[t'| t]$ is  $\theta$ -free; so that,  $E_{\theta}[t'|t]$  is a statistic. Now, t' =  $E_{\theta}[t'|t] + \varepsilon_{t't}$ , where  $\varepsilon_{t',t}$  is the error due to the regression of t' on t and  $E_{\theta}[t'|t] = h(t) = t''$ is the regression function of t 'on t. Also,  $E_{\theta}[h(t)] = E_{\theta}[t''] = g(\theta)$ , say.

Furthermore, since h(t) and  $\varepsilon_{t,t}$  are uncorrelated (this is immediate from regression theory), we have

Var(t') = Var(t'') + Var(t'-h(t)),

hence,

Var (t')  $\leq$  Var (t")

with equality holding if and only if t' = h(t). (i.e. t' and

t are functionally dependent) with probability one.

We note that the Rao-Blackwell theorem gives us a way of improving an arbitrary unbiased estimator t' of  $g(\theta)$  when a sufficient statistic t for  $\theta$  exists and is functionally independent of  $\theta$ ; t" calculated from t' and t by the above method is unbiased for the parameter  $\theta$  and it has variance less than the variance of t', for all sufficient statistics t. We also note that in seeking a uniformly minimum variance unbiased estimator of  $g(\theta)$ , we only need consider unbiased estimators which are functions of the sufficient statistic t, since any other estimator for  $\theta$  cannot be a uniformly minimum variance unbiased estimator. The above idea will now be used in the second theorem.

## 2.7.4 The Lehmann-Scheffé Theorem.

We have previously mentioned that completeness confers the uniqueness property on sufficient statistics. The following theorem, due to Lehmann and Scheffé, characterizes this uniqueness property.

<u>Theorem</u>. If the statistic t is complete and sufficient for a parameter  $\theta_{j}$  and  $g(\theta)$  is an estimable function of  $\theta$ , then the function  $g(\theta)$  has a uniformly minimum variance unbiased estimator h(t), a function of t only, which can be determined essentially in either of the following two ways.

(i) If t' is any unbiased estimator of the function  $g(\theta)$ , which is functionally independent of the statistic t, then

 $h(t) \equiv E_{\theta}[t'|t]$ 

(ii) h(t) is the solution of the equation

 $E_{\theta}[h(t)] = g(\theta)$ .

<u>Proof.</u> Define  $U = U_t \bigcup \overline{U}_t$ , where  $U_t$  is the set of all unbiased estimators of  $g(\theta)$  which are functionally dependent on t and  $\overline{U}_t$ is the set of all unbiased estimators of  $g(\theta)$  which are functionally independent of t. U is the set of unbiased estimators. Then,  $U_t \cap \overline{U}_t = \Phi$ ,  $\phi$  denotes the null set. Since g(.) is an estimable function and  $U \neq \phi$ ,  $U_t$  and  $\overline{U}_t$  cannot both be empty. (i) Assume  $\overline{U}_t$  is not empty and let  $t' \in \overline{U}_t$ , then, by the Rao-Blackwell theorem,

$$t'' = E_{\theta} \left[ t' | t \right] = h(t), \text{ say,}$$

is an unbiased estimator for  $g(\theta)$ , and h(t) is an element of  $U_t$  which has uniformly smaller variance than t'. Moreover, since t is complete for  $\theta$ , h(t) is an essentially unique member of  $U_t$  and is determined by the relationship

$$E_{\theta}[h(t)] = g(\theta)$$
.

Thus h(t) is a uniformly minimum variance unbiased estimator of  $g(\theta)$ .

(ii) Suppose  $U_t$  is not empty, then as in (i)  $U_t$  contains essentially only one element t" = h(t), which is determined by

$$E_{\theta}\left[h(t)\right] = g(\theta)$$

and hence t" is a uniformly minimum variance unbiased estimator for  $g(\theta)$ . We are now in a position to prove the

uniqueness of a uniformly minimum variance unbiased estimator for a parametric function.

<u>Proof</u>: Suppose there are two such estimators of  $g(\theta)$ ,  $t_1$ and  $t_2$ , say. Then  $t_3 = (t_1 + t_2)/2$  is also unbiased for  $g(\theta)$ .

Let V by the common variance of  $t_1$  and  $t_2$  Then,

$$Var(t_3) = (2V + 2 Cov (t_1, t_2))/4$$
  
= (V + Cov(t\_1, t\_2))/2 .

Now, by the Cauchy-Schwarz inequality

Cov (t<sub>1</sub>, t<sub>2</sub>)  $\leq \sqrt{(Var t_1 Var t_2)} = V$ .

Thus  $Var(t_3) \leq V$ , which contradicts the assumption that t1 and t2 are minimum variance unbiased estimators. Hence,  $Cov(t_1, t_2) = V$ 

and 
$$Var(t_3) = V$$
.

This can only be so if

$$t_1 - g(\theta) = k(\theta) (t_2 - g(\theta)) ,$$

i.e., if t1 is proportional to t2, but then,

Cov  $(t_1, t_2) = k(\theta) \cdot V = V$ 

hence k ( $\theta$ ) = 1, and so t<sub>1</sub> = t<sub>2</sub>.

Q.E.D.

We now give an example of how to determine the uniformly minimum variance estimator.

Example: Let  $X_1, X_2, \ldots, X_n$  be a random sample from a uniform distribution on  $(0, \theta)$ ,  $0 < \theta < \infty$ , let  $t = \max(X_1, X_2, \ldots, X_n)$ . It can be shown that t is complete and sufficient for the parameter  $\theta$ .

We are to determine a uniformly minimum variance unbiased estimator of  $\theta^r$ ,  $r \neq -n$ .

By the Lehmann-Scheffé theorem some function of the statistic t will be such an estimator. Assume t<sup>r</sup> is an estimator of  $\theta^r$ . then,

$E_{\theta}[t^{r}]$	=	$\int_{0}^{\theta} t^{r}g(t)dt$
	н	$\int_{0}^{\theta} t^{r} \cdot \underline{nt}^{n-1} dt$
	=	$\frac{1}{\theta^{n}} \left[ \frac{n}{n-r} t^{n-r} \right]_{0}^{\theta}$
	=	$\frac{n}{n-r} \theta^r$ .

Thus  $\frac{n-r}{n} t^r$  is unbiased for  $\theta^r$  and will be the uniformly

minimum variance unbiased estimator.

### 2.8 Efficiency

We noticed earlier that the sample mean  $\bar{\mathbf{x}}$  and median  $\mathbf{x}_1$  were two consistent estimators for the population mean of a normal distribution with mean  $\mu$  and variance  $\sigma^2$ , and that these two estimators were both asymptotically normal with the same mean  $\mu$  and variance  $\sigma^2/n$  and  $\pi\sigma^2/2n$  respectively.

There is a large and important class of consistent estimators whose sampling distributions are asymptotically normal and for which the variance decreases as the sample size increases, so that the properties of such estimators can be characterized by the variance and bias. The bias will not be considered, since it must tend to zero as the sample size increases, or in the first case we could have made a correction for the bias, as noted earlier. We are not so fortunate when it comes to the variance, since the variance is invariant under a constant shift. We therefore seek a criterion which will enable us to choose an estimator from among competing ones.

The concept of efficiency as defined by R. A. Fisher (1921) requires that the fixed value to which n times the estimator tends, shall be as small as possible. An efficient estimator will be such that this property is satisfied.

If the variance of any efficient estimator is known, then, the efficiency of any other estimator of the same parameter is calculated as the ratio of the variance of the efficient estimator to that of the other estimator. <u>Definition</u>: In the class of all unbiased estimators of a parameter  $\theta$ ,  $\hat{\theta}$  is efficient if

$$\operatorname{Var}(\hat{\theta}) \leq \operatorname{Var}(\hat{\theta'})$$

that is,  $\theta$  is efficient if it is a uniformly minimum variance unbiased estimator.

Example. Consider again a normal population with mean  $\mu$  and variance  $\sigma^2$ . Let  $\mathbf{x}_1, \mathbf{x}_2, \dots \mathbf{x}_n$  be a random sample of size n chosen from this population. As mentioned earlier, the sample mean  $\bar{\mathbf{x}}$  and the sample median  $\mathbf{x}_1$  are both unbiased consistent estimators for the population mean  $\mu$ . However, the variance of the sample mean is  $\sigma^2/n$  and that of the sample median is  $\pi\sigma^2/2n$ . Hence the efficiency of the median is

$$\frac{\sigma^2/n}{\pi\sigma^2/2n} = \frac{2}{\pi} \approx .637.$$

This means that from 1000 observations the sample median will produce a value which will be as accurate as that produced by the sample mean for only 637 observations, and therefore as mentioned in the discussion of consistency, the sample mean would be a better choice as an estimator for the population mean in the normal population.

Sometimes, it is more difficult to calculate an efficient estimator than an inefficient one, and the labour involved in calculating this efficient estimator might outweigh the cost of taking additional observations needed for the inefficient estimator. In such cases, the additional observations may be taken and the inefficient estimator chosen as the estimator for the population parameter.

### 2.8.1. Properties of efficient estimators.

(i) The correlation between any two statistics, both efficient estimators of the same parameter  $\theta$ , tends to + 1, as the sample size increases. For if  $t_1$  and  $t_2$  are two such statistics whose variances are  $\sigma^2/n$ , say, and the correlation is r, we choose  $t_3 = (t_1 + t_2)/2$ . Then, the statistic  $t_3$ will be an efficient estimator of the parameter  $\theta$ . The variance of  $t_3$  is  $(1 + r)\sigma^2/n$ , , and this, by hypothesis, must be at least  $\sigma^2/n$ . Hence, r must be at least + 1, but r is at most + 1, therefore r = + 1. This property implies that for large samples all efficient statistics (estimators) are equivalent.

(ii) If we relax the constraint of unbiasedness in the definition of efficiency and replace the variance by the mean square error, we note that since maximum likelihood estimators have asymptotically smaller variances (see Kendall and Stuart, Vol.II) they are therefore efficient.

#### CHAPTER III

#### MAXIMUM LIKELIHOOD ESTIMATION

3.1 <u>Introduction</u>. One of the most important methods of estimation of population parameters is the method of Maximum Likelihood. This method of estimation was formally introduced by Fisher, R.A. in a short paper in 1912, and was further developed in a series of papers (Fisher (1922, 1925, 1934)) by the same author. Several other authors have since made important contributions to the theory.

Consider a random sample  $X = (x_1, x_2, ..., x_n)'$  taken from a population with parameter  $\theta$ , and known probability density (mass) function (p.d.f.) f. The joint probability of the observations regarded as a function of the unknown parameter  $\theta$  is called the likelihood function, and is denoted by  $L(X \mid \theta)$ , or simply by  $L(\theta)$ , and is given by

$$L(\theta) = \frac{n}{n} f(x_i \mid \theta), i=1,2,\ldots,n.$$

We note here that since L  $(\theta)$  is the product of p.d.f.'s (p.m.f.'s) it is necessarily positive.

The method of Maximum Likelihood in its simplest form leads us to take as estimate of the parameter  $\theta$ , that value  $\hat{\theta}$ which lies within the range of admissible values of the parameter  $\hat{\theta}$  and makes the likelihood function  $L(\theta)$  as large as possible,
that is,

 $L(\theta) \ge L(\theta)$  for all values of  $\theta \in \Theta$ 

the parameter space, and  $\theta \in \Theta$ . If we assume that  $L(\theta)$  is a twice differentiable function (in the range of  $\theta$ ) with respect to  $\theta$ , then under suitable conditions we can use calculus to find  $\hat{\theta}$ . Namely, we solve the following constrained system

$$\frac{\partial}{\partial \theta} L(\theta) = 0,$$
 (3.1.1)

subject to

$$\frac{\partial^2}{\partial \theta \partial \theta} L(\hat{\theta}) < 0^*.$$
 (3.1.2)

We note that the condition (3.1.2) is only a sufficient condition for maximality of  $\hat{\theta}$ , and that it is not a necessary condition. For consider

$$g(\theta) = -\theta_1^4 - \theta_2^4$$
,  $\theta = (\theta_1, \theta_2)'$ ,

then

$$\frac{\partial}{\partial \theta} g(\theta) = (-4\theta_1^3, - 4\theta_2^3)$$

and

$$\frac{\partial}{\partial \theta} \frac{\partial}{\partial \theta} g(\theta) = -12 \begin{bmatrix} \theta_1^2 & 0 \\ 0 & \theta_2^2 \end{bmatrix}$$

 $\theta = (0, 0)'$  can be shown to be the point of maximum of  $g(\theta)$ ,

however

$$\frac{\partial}{\partial \theta}$$
,  $g(\theta) |_{\theta} = (0,0)'$ 

The matrix inequality < indicates negative definiteness.

is not negative definite.

In many situations the system of equations obtained by setting  $\frac{\partial L(\theta)}{\partial \theta} = 0$  is quite complicated and rather difficult to solve. In some of these cases it is much easier to work with the logarithm of the likelihood function denoted by  $\ell(\theta)$ . This is possible because a positive function and its logarithm attain their maxima (minima) at the same point. Thus we solve the following constrained system;

$$\frac{\partial}{\partial \theta} \ell(\theta) = 0,$$
 (3.1.3)

subject to

We know from (1.7) that

$$\frac{\partial}{\partial \theta} \& (\theta) = S,$$

and from (1.8) that

$$- E\left[\frac{\partial^2}{\partial \theta \partial \theta}, \ell(\theta)\right] = J(\theta).$$

3.2 Regularity Assumptions

Before discussing iterative methods for the solution of the likelihood equation, we mention here some regularity assumptions which we place on the likelihood function  $L(\theta)$ , and discuss the properties of the score vector,  $S(\theta)$ , and the information matrix,  $J(\theta)$ . RA1.  $\theta \in \bigoplus_{\ell}$ , where  $\bigoplus_{\ell}$  is a subset of  $|\mathbb{R}^{\ell}$  (Euclidean *l*-space) possibly equal to  $|\mathbb{R}^{\ell}$ .

RA2. For all samples X, 
$$\frac{\partial L(\theta)}{\partial \theta_{j}}$$
; j = 1,...,l, exists

for all  $\theta \in \Theta_{\ell}$  except perhaps for a set of points in  $\Theta_{\ell}$  with probability measure zero.

- RA3.  $\int L(\theta) dx$  (or  $\Sigma L(\theta)$ ) can be differentiated under the X integral (summation) sign with respect to  $\theta_j$ , j = 1, ..., l.
- RA4.  $\int t_i L(\theta) dx$  (or  $\Sigma t_i L(\theta)$ ) can be differentiated under X integral (summation) sign with respect to  $\theta_j$ , j = 1, ..., l, where  $(t_i)$ , i=1, ..., m is a sufficient statistic for the parameter  $\tau(\theta)$ .

RA5. The score vector  $S(\theta) = \frac{\partial}{\partial \Theta} \ell(\Theta)$  has a positive definite covariance matrix for each  $\Theta \epsilon \Theta$ .

3.3. Properties of the Score Vector

(1) Let  $x_i (m_i \times 1) = 1, \dots, q$  be independent vectors, where the distribution of each  $x_i$  involves the same  $\theta$ . Let  $X' = (x'_1, x'_2, \dots, x'_q)$ . Denote the score vectors of  $x_i$ and X by  $S_i(\theta)$  and  $S(\theta)$  respectively. Then

 $S(\theta) = \sum_{i=1}^{q} S_{i}(\theta) . \qquad (3.3.1)$ <u>Proof</u>:  $L(\theta, X) = \frac{q}{i=1} L(\theta, x_{i}).$ Hence,  $\log L(\theta, X) = \sum_{i=1}^{q} \log L(\theta, x_{i}).$ 

Whence, 
$$\frac{\partial}{\partial \theta_{j}}$$
 log  $L(\theta, X) = \sum_{i=1}^{q} \frac{\partial}{\partial \theta_{j}}$  log  $L(\theta, x_{i})$ ,  $j = 1, \dots, l$ .

Thus

$$s(\theta) = \sum_{i=1}^{q} s_i (\theta).$$

(ii) For all values of  $\theta \in \Theta_{\ell}$ , the expected value of the score vector is zero.

<u>Proof</u>: Since  $L(\theta)$  is a probability density, we have

 $\int L(\theta) dx = 1$ . X Hence, using RA2 and RA3, we see that

But,

$$\int_{X} \frac{\partial}{\partial \theta_{j}} L(\theta) dx = \int_{X} \frac{1}{L(\theta)} \frac{\partial L}{\partial \theta_{j}} \frac{\partial L(\theta)}{\partial \theta_{j}} L(\theta) dx$$
$$= \int_{X} \int_{y} (\theta) L(\theta) dx .$$

 $\int_{\mathbf{X}} \frac{\partial}{\partial \Omega} L(\theta) d\mathbf{x} = 0, j = 1, \dots, \ell.$ 

Hence  $E[S(\theta)] = 0$ .

(iii) The covariance matrix  $\Sigma_{t,S}$  between S and any unbiased estimator t is given by

$$\Sigma_{t,S} = \frac{\partial}{\partial \theta} \tau(\theta),$$

where  $\tau(\theta) = E[t]$ . In particular, for  $\tau(\theta) = \theta$ , we have  $\Sigma_{t,S} = I$ .

Proof.

$$E\left[t\right] = \int t \cdot L(\theta) dx = \tau(\theta) dx$$

Thus using RA3 and RA4, we have

$$\int_{X} t \frac{\partial}{\partial \theta} L(\theta) dx = \frac{\partial}{\partial \theta} \tau(\theta),$$

i.e., 
$$\int_{X} t \frac{1}{L(\theta)} \frac{\partial}{\partial \theta} L(\theta) L(\theta) dx = \frac{\partial}{\partial \theta} \tau(\theta) ,$$
  
or 
$$\int_{X} tS(\theta) L(\theta) dx = \frac{\partial}{\partial \theta} \tau(\theta)$$
$$= \sum_{t,S} .$$

# 3.4 Properties of the Information Matrix

(i) Additivity:

If  $J_1(\theta)$ ,  $J_2(\theta)$ , ...,  $J_n(\theta)$  are the information matrices for n independent random experiments, then the information matrix  $J(\theta)$  for the combined experiment is given by

$$J(\theta) = \sum_{i=1}^{n} J_i(\theta)$$
.

If  $X' = (Z'_1, Z'_2, ..., Z'_n)$ , where X is the ordered observations from the combined experiment, and  $Z'_1, Z'_2, ..., Z'_n$  is a random sample from a distribution with p.d.f. (p.m.f.)  $f(x, \theta)$ , then

$$J(\theta) = nJ_{0}(\theta),$$
  
where  $J_{0}(\theta) = E \left[\frac{\partial}{\partial \theta} \log f \left[\frac{\partial}{\partial \theta} \log f\right]'\right]$ 

 $J_{0}^{}\left( \theta\right) \,$  is called the information matrix per unit observation. Proof

Let  $S_i(\theta)$  i = 1,...,n be the score vector for the ith experiment, and let  $S(\theta)$  be the score vector for the combined experiment.

Then 
$$J(\theta) = E[S(\theta) S(\theta)']$$
  

$$= E\begin{bmatrix} n & S_{i}(\theta) & n & S_{j}(\theta)' \\ i=1 & S_{i}(\theta) & j=1 & J_{j}(\theta)' \end{bmatrix}$$

$$= E\begin{bmatrix} n & S_{i}(\theta) & S_{i}(\theta)' + 2 & n & S_{i}(\theta) & S_{j}(\theta)' \end{bmatrix}$$

$$= n & E\begin{bmatrix} S_{i}(\theta) & S_{i}(\theta)' + 2 & j & S_{i}(\theta) & S_{j}(\theta)' \end{bmatrix}$$

$$= n & E\begin{bmatrix} S_{i}(\theta) & S_{i}(\theta)' & J & J_{i}(\theta) & J & J_{i}(\theta) & J & J_{i}(\theta) & J & J_{i}(\theta) & J_{i$$

independent for  $i \neq j$ , being functions of  $Z_i$  and  $Z_j$  respectively, with  $Z_i$  and  $Z_j$  independent by hypothesis.

(ii) If in addition to the assumptions RA1 to RA5 we assume that  $\frac{\partial^2}{\partial \Theta \partial \Theta} L(\Theta)$  exists almost everywhere (x) for all  $\Theta \in \bigoplus_{\ell}$  and if  $\int_X L(\Theta) dx$  can be differentiated twice under the integral sign, then

$$J(\Theta) = E\left[-\frac{\partial^2}{\partial \theta \partial \theta}, L(\Theta)\right].$$

$$\frac{\text{Proof:}}{\partial \theta} \quad \frac{\partial}{\partial \theta} \log L(\theta) = \frac{1}{L(\theta)} \frac{\partial}{\partial \theta} L(\theta) ,$$
thus,  $\frac{\partial^2}{\partial \theta \partial \theta} \log L(\theta) = \frac{1}{L(\theta)} \frac{\partial^2}{\partial \theta \partial \theta} L(\theta) - \frac{1}{L^2(\theta)} \frac{\partial}{\partial \theta} L(\theta) \left[ \frac{\partial}{\partial \theta} L(\theta) \right]'$ 
and  $E\left[ - \frac{\partial^2}{\partial \theta \partial \theta} \log L(\theta) \right] = E\left[ \frac{1}{L^2(\theta)} \frac{\partial}{\partial \theta} L(\theta) \left[ \frac{\partial}{\partial \theta} L(\theta) \right]' \right]$ 

$$-E\left[ \frac{1}{L(\theta)} \frac{\partial^2}{\partial \theta \partial \theta} \log L(\theta) \right]$$

if both expectations on the right hand side exist.

But, 
$$E\left[\frac{1}{L^{2}(\theta)}\frac{\partial}{\partial \theta}L(\theta)\left[\frac{\partial}{\partial \theta}L(\theta)\right]'\right] = \int_{X} \frac{1}{L^{2}(\theta)}\frac{\partial}{\partial \theta}L(\theta)\left[\frac{\partial}{\partial \theta}L(\theta)\right]'L(\theta)dx$$
  
$$= \int_{X} \frac{1}{L(\theta)}\frac{\partial}{\partial \theta}L(\theta)\frac{1}{L(\theta)}\left[\frac{\partial}{\partial \theta}L(\theta)\right]'L(\theta)dx$$
$$= \int_{X} \frac{\partial}{\partial \theta}L(\theta)\left[\frac{\partial}{\partial \theta}L(\theta)\right]'L(\theta)dx$$
$$= E\left[S(\theta)S(\theta)'\right]$$
$$= J(\theta)$$

41

.

.

and

$$E\left[\frac{1}{L(\theta)} \begin{array}{c} \frac{\partial}{\partial \theta \partial \theta} L(\theta)\right] = \int_{X} \frac{1}{L(\theta)} \frac{\partial}{\partial \theta \partial \theta} L(\theta) L(\theta) dx$$
$$= \int_{X} \frac{\partial}{\partial \theta \partial \theta} L(\theta) dx$$
$$= \frac{\partial}{\partial \theta \partial \theta} \int_{X} L(\theta) dx$$
$$= \frac{\partial}{\partial \theta \partial \theta} \int_{X} L(\theta) dx$$
$$= \frac{\partial}{\partial \theta \partial \theta} \int_{X} L(\theta) dx$$

= O.

.

,

# 3.5 <u>Iterative Methods for the Solution of the Maximum Likeli-</u> hood Equation: The Single Parameter Case.

In many instances it is very difficult to find an explicit solution of the Maximum Likelihood equation. In these cases the Maximum Likelihood Estimate is approximated by using iterative techniques. The most popular technique for approximating the Maximum Likelihood Estimate is the method of scoring for parameters. This technique was suggested by Fisher (1925). Several other methods are also used. These include the Newton-Raphson method, the fixed derivative Newton method and the Regula Falsi method and its modifications.

Kale (1961) discusses the Newton-Raphson method, the fixed derivative Newton method and the method of scoring for parameters. He gives sufficient conditions for the convergence of these iterative schemes and applies the three techniques to the solution of a particular problem. Barnett (1966) also discusses these three methods together with the Regula Falsi method, and notes that the Regula Falsi method will be preferred to these three schemes in cases where the likelihood equation has multiple roots. This is so since the Regula Falsi method can be made to scan the range of the likelihood equation to find all the relative maxima.

We now give a brief discussion of a general iterative scheme of which the Newton-Raphson method, the method of scoring

for parameters, the fixed derivative Newton method and the Regula Falsi method are special cases.

The idea is to choose an initial solution  $t_0$ , which lies in a neighbourhood of the Maximum Likelihood Estimate, and use an iterative scheme to find the Maximum Likelihood Estimate.  $t_0$  is usually chosen to be an unbiased, consistent estimator.

A possible iterative scheme is as follows:

Define

$$\phi (\theta) = \theta - \psi (\theta) \underline{d} \ell(\theta),$$

where  $l(\theta) = \log L(\theta)$ .

Consider the iterative scheme

$$t_{n+1} = \phi(t_n) \qquad (3.5.1)$$
$$= t_n - \psi(t_n) \left[ \frac{d}{d\theta} \ell(\theta) \right]_{\theta} = t_n \qquad (3.5.2)$$

Let  $e_n = |t_n - \hat{\theta}|$  be the error at the nth iteration, then we choose  $\psi$  ( $\theta$ ) such that  $e_{n+1} < e_n$  and  $e_n \rightarrow 0$  as  $n > \infty$ . Householder (1953) has shown that it is sufficient that: (1) There exists an  $\varepsilon$ -neighbourhood N ( $\hat{\theta}$ ) of  $\hat{\theta}$ , such that if  $\theta_1$  and  $\theta_2$  are in N ( $\hat{\theta}$ ), then for some  $k \ge 0$ , we have

$$|\phi (\theta_1) - \phi (\theta_2)| \leq k < 1,$$

(2)  $t_1 \in \mathbb{N}_{\epsilon}(\theta)$ 

for (3.5.2) to converge.

The Newton-Raphson method, the method of scoring for parameters, the fixed derivative Newton method and the secant

method\* are all special forms of (3.5.1). In the Newton-

thod, 
$$\psi(\theta) = \frac{d^2}{d\theta^2} \ell(\theta)$$
.

Raphson me

In the method of scoring for parameters  $\psi(\theta) = E_{\theta} \left[ \frac{d^2}{d\theta^2} \ell(\theta) \right]$ , i.e.,  $\psi(\theta) = -J(\hat{\theta})$ . In the fixed derivative Newton method  $\psi(\theta) = k$ , where k is a constant such that the conditions for convergence of (3.5.2) hold.

 $\psi(\theta) = \left[\frac{d}{d\theta}\ell(\theta) - \left(\frac{d}{d\theta}\ell(\theta)\right)_{\theta} = t\right] / (\theta - t) \quad \text{in the secant}$ method, here t is a consistent estimator of  $\hat{\theta}$ .

It is well known that the Newton-Raphson method converges quadratically, i.e.,  $e_{n+1} \, \,^{\alpha} \, e_n^2$ , whereas the orders of convergence of the other methods lie in the half open interval [1,2). Thus it might seem that the Newton-Raphson method should be preferred to the other methods. However, as pointed out by Barnett (1966) the first three methods listed above all have the same undesirable qualities, and are therefore not recommended in cases where it is suspected that the likelihood equation has multiple roots. Barnett (1966) shows that the Regula Falsi method (and hence its modified forms) works well in many cases where the first three methods fail. For further discussions of iterative techniques the reader is referred to Conte and deBoor (1972) and Hildebrand (1956).

We now use the Newton-Raphson method, the Modified

\* This is a modified form of the Regula Falsi method.

Regula Falsi method and the Secant method to solve the genetical example given by Fisher (1954) based on Carver's data for two factors in corn; Starchy vs. Sugary and Green vs. White.

We must evaluate the Maximum Likelihood Estimate  $\hat{\theta}$  of  $\theta$  when the probability of belonging to one of the four classes listed below are

Starchy

# Sugary

 $\frac{\text{Green}}{p_1 = (2+\theta)/4} \quad \frac{\text{White}}{p_2 = (1-\theta)/4} \quad \frac{\text{Green}}{p_3 = (1-\theta)/4} \quad p_4 = \theta/4 .$ 

A random sample of size 3839 was taken and the numbers falling into each class are given in the following table:

Star	chy	Sugary		
Green	White	Green	White	
= 1997	b = 906	c = 904	d = 32	

Since the distribution is multinominal, we have

$$L(a,b,c,d|\theta) = k(2 + \theta)^{a}(1 - \theta)^{b} + c^{b}\theta^{d}$$
,

where k is a constant.

Hence

а

$$\frac{\mathrm{d}}{\mathrm{d}\theta} \quad \frac{\ell(\theta)}{2+\theta} = \frac{\mathrm{a}}{1-\theta} - \frac{\mathrm{b}}{\theta} + \frac{\mathrm{c}}{\theta} - \frac{\mathrm{d}}{\theta} \quad (3.5.3)$$

Assuming that  $0 < \theta < 1$ , it is easy to show that the Maximum

Likelihood Estimate  $\hat{\theta}$  can be obtained by solving the following quadratic equation;

 $n\theta^2 + (2b+2c+d-a)\theta-2d = 0.$ 

 $\hat{\theta}$  is taken to be the solution which lies in  $\bigoplus_{\ell}$  and maximizes  $\ell(\theta)$ . For the above data Fisher (1950) has shown that  $\hat{\theta} = .035712$  to six decimal places. Norton (1956) and Kale (1961) use the Newton-Raphson method starting from  $t_0 = (a-b-c+d)/n = .057046$ , and after five iterations obtain  $\hat{\theta} = .035712$  and .035713 respectively. We shall start iterating from  $t_0 = (a+5d-b-c)/2n$ = .045194. In the modified Regula Falsi and secant methods we shall use as our two initial estimates  $t_{-1} = 4d/n = .033342$  and  $t_0 = (a+5d-b-c)/2n = .045194$ . It is observed that  $t_{-1}$  and  $t_0$  are both unbiased and consistent estimators for  $\theta$ , and that the choice of  $t_{-1}$  is justified since

 $\begin{bmatrix} \underline{d} \ \ell(\theta) \end{bmatrix}_{\theta = t_{-1}} = 69.446846 \text{ and } \begin{bmatrix} \underline{d} \ \ell(\theta) \\ \overline{d\theta} \end{bmatrix}_{\theta = t_{0}} = -211.178969 ,$ i.e.,  $\hat{\theta} \boldsymbol{\epsilon} \begin{bmatrix} t_{-1}, t_{0} \end{bmatrix}$ . The results for these methods are presented in the table below.

Iterates	I	II	III .
t_1	_	.033342	.03334.2
to	.045194	.045194	.045194
t <sub>1</sub>	.033546	.036275	.036275
t <sub>2</sub>	.035591	.035379	.035580
t <sub>3</sub>	.035712	.035717	.035714
t4	_	.035712	.035712

Successive Iterations by the Three Different Methods.

I denotes the Newton-Raphson Method.

II denotes the modified Regula Falsi Method.

III denotes the secant method.

It is observed that the secant method performs just as well as the Newton-Raphson method. We also observe that the Newton-Raphson method with our choice of  $t_0$  takes one less iteration than with the  $t_0$  used by Kale (1961) and Norton (1956).

We note here that the Newton-Raphson method, the method of scoring for parameters, and the fixed derivative Newton method can be applied to the multiparameter case.

# 3.6 Properties of MLE's Obtained from Dependent Observetions

Under certain regularity conditions (mentioned earlier) Cramér (1946) proved that the MLE obtained from independent observations is:

(i) weakly consistent, and

(ii) asymptotically efficient and normally distributed. These results were later extended by other authors, some of whom relaxed a few of Cramér's condition, while others imposed additional conditions on the likelihood equation.

In recent times Silvey (1961), Bar-Shalom (1970) and Bhat (1973) considered the asymptotic properties of MLE'S obtained from dependent observations. Bar Shalom (1970) imposed two conditions on the likelihood equation in addition to those stated by Cramér. He then used a version the the weak law of large numbers for dependent random variables to prove the convergence of  $b_0$  and  $b_1$  (defined below). Bhat (1973) used the central limit theorem for Martingale sequences to prove the convergence of  $b_0$  and  $b_1$ .

We now list the Regularity conditions as given by Bar-Shalom (1970).

# Notation and Listing of Regularity Conditions.

Let the set of (possibly dependent) observations be  $Z^n = (z_1, z_2, ..., z_n)$ , where  $z_i$ , i = 1, ..., n, are real valued random variables with joint probability density function (p.d.f) With respect to a  $\sigma$ -finite (i.e., finite with respect to the  $\sigma$ -algebra generated by  $Z^n$ ) measure  $\mu_n$ , given by

$$P(z_1, z_2, \dots, z_n | \theta) = P(Z^n | \theta)$$
(3.6.1)

where  $\theta$  is a real valued constant with unknown, but true

value  $\theta_0 \in \Theta$ , with  $\Theta$  an open subset of  $|\mathbb{R}^1$ . The Borel measurable function  $\hat{\theta}_n(Z^n)$  obtained by maximization of the likelihood function

$$L_{n}(\theta) = P(Z^{n}|\theta) = \prod_{i=1}^{n} P_{i}(\theta)$$
,

where  $p_i(\theta) = L_i(\theta) / L_{i-1}(\theta)$  is used as an estimator for  $\theta$ . We now list the following regularity conditions which are needed to prove the asymptotic properties of the MLE.

RC1. For almost all 
$$Z^{k}$$
,  $\frac{\partial i}{\partial \theta i}$  log  $p_{k}(\theta)$ , i=1,2,3, exist

for all  $\theta \in \Theta$ .

RC2. 
$$E\left[\frac{\partial}{\partial\theta} \log p_{k}(\theta)\right] = 0$$
.  
RC3.  $J_{k}(\theta) = E\left[\frac{\partial}{\partial\theta} \log p_{k}(\theta)\right]^{2}L_{k}(\theta) \prod_{i=1}^{k} \mu(dy_{i}) \leq C_{i} < \infty$ ,  
where  $C_{1}$  is independent of  $\theta$ , and  $J_{k}(\theta)$  is Fisher's  
information measure (see for example Rao (1973)).  
RA4.  $E\left[\frac{\partial}{\partial\theta^{2}} \log p_{k}(\theta)\right] = -J_{k}(\theta_{0})$ .  
RC5. There exists a function,  $B_{k}(Z^{k})$ , measurable with respect  
to the product measure  $\prod_{i=1}^{k} \mu(dy_{i})$ ,  
such that  $\left|\frac{\partial}{\partial\theta^{3}} \log p_{k}(\theta)\right| \leq B_{k}(Z^{k})$  for all  $\theta$ ,  
and  $B_{k}(Z^{k})$  is finite almost everywhere, i.e., there  
exists a constant,  $N < \infty$ , independent of  $\theta$  and k,

such that for all  $\varepsilon > 0$ ,  $P(B_k(Z^k) > N) < \varepsilon$ .

RC6.  $E\left[\frac{\partial}{\partial \theta} \log p_{i}(\theta) \begin{array}{c} \partial \\ \partial \\ \theta \end{array} \right] \log p_{j}(\theta) = 0, \text{ for all } i \neq j.$ A milder version of RC6 is;

RC7. 
$$\lim_{|i-j| \to \infty} E \left[ \frac{\partial}{\partial \theta} \log p_i(\theta) \frac{\partial}{\partial \theta'} \log p_j(\theta) \right] = 0$$
  
RC8. 
$$\operatorname{Var} \left[ \frac{\partial}{\partial \theta^2} \log p_i(\theta) \right] \leq C_2 < \infty$$

where  $C_2$  is a constant independent of i, and

$$\lim_{|\mathbf{i}-\mathbf{j}|\to\infty} \operatorname{Cov}\left[\frac{\partial}{\partial\theta^2} \log p_{\mathbf{i}}(\theta) \frac{\partial}{\partial\theta^2} \log p_{\mathbf{j}}(\theta)\right] = 0.$$

<u>Note</u>: In the above regularity conditions all the expectations are evaluated at,  $\theta = \theta_0$ , the true value. Conditions RC1 to RC5 are similar to the conditions RA1 to RA5 given by Cramér (1946). Conditions RC6 to RC8 are needed because of the possible dependence among the observations.

We are now in a position to state and prove the following theorem.

<u>Theorem</u>: Under regularity conditions RC1 to RC8, the Maximum Likelihood Estimator  $\hat{\theta}_n(Z^n)$  is weakly consistent, i.e.,

 $\hat{\theta}_n(z^n) \xrightarrow{p} \theta_o,$ 

(where  $\underline{P}$  denotes convergence in probability) and asymptotically efficient, i.e., for sufficiently large n the mean square of  $\hat{\theta}_n(Z^n)$  is equal to the reciprocal of the total information, i.e.,

$$\mathbb{E}\left[\left(\hat{\theta}_{n}-\theta_{0}\right)^{2}\right] = \begin{bmatrix} n \\ \Sigma \\ i=1 \end{bmatrix} J_{i}(\theta_{0}) = 1.$$

Proof: Following Cramér, by conditions RC1 to RC6

$$\frac{\partial}{\partial \theta} \log p_{k}(\theta) = \left[ \frac{\partial}{\partial \theta} \log p_{k}(\theta) \right]_{\theta=\theta_{0}}^{+} (\theta-\theta_{0}) \left[ \frac{\partial}{\partial \theta^{2}} \log p_{k}(\theta) \right]_{\theta=\theta_{0}}^{+} \frac{\lambda}{2} (\theta,\theta_{0}) \cdot (\theta-\theta_{0})^{2} B_{k}(z^{k}) , \quad (3.6.2)$$

where  $B_{\mu}$  is defined above, and  $|\lambda| < 1$ . Thus the likelihood equation (3.6.2) may, after multiplication by  $\frac{1}{n}$ and summing, be written as

 $\frac{1}{n} \frac{\partial}{\partial \theta} L_n(\theta) = b_0 + (\theta - \theta_0) b_1 + \frac{\lambda}{2} (\theta - \theta_0)^2 b_2 = 0 ,$ where  $b_0 = \frac{1}{n} \sum_{i=1}^{n} \left[ \frac{\partial}{\partial \theta} \log p_i(\theta) \right]_{\theta = \theta_0}$  $b_{1} = \frac{1}{n} \sum_{i=1}^{n} \left[ \frac{\partial^{2}}{\partial \theta^{2}} \log p_{i}(\theta) \right]_{\theta = \theta_{0}}$  $b_2 = \frac{1}{n} \sum_{i=1}^{n} B_i(Z^i).$ 

We now use a version of the weak law of large numbers for dependent random variables (Parzen (1969), p. 418) to prove the convergence of  $b_0$  and  $b_1$ . This law states that if for all i the random variable x<sub>i</sub> has finite mean and bounded variance  $\lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} E(x_i) < \infty$ , then a sufficient and

 $\frac{1}{n} \sum_{i=1}^{n} x_i - \frac{1}{n} \sum_{i=1}^{n} E(x_i) \xrightarrow{p} 0,$ condition for

 $\lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{\Sigma} Cov (x_i, x_j) = 0.$ is

Using this theorem, it follows from RC2, RC3 and RC7 that, b<sub>0</sub>  $\xrightarrow{P}$  0,

and from RC4 and RC8 that,  $b_1 \xrightarrow{p} - J(\theta_0)$ , where

 $J(\theta_0) = \lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} J_i(\theta_0) \text{ exists by virtue of}$ 

RC3. From RC5, it follows that for all  $\varepsilon > 0$ , there exists N such that

We now use these results to show that the estimator  $\boldsymbol{\hat{\theta}}_n(\boldsymbol{Z}^n)$  is weakly consistent.

Let  $\delta$  and  $\varepsilon$  be given arbitrarily small positive numbers. For sufficiently large  $n(n > n_0(\delta, \varepsilon), say)$ 

we have .

$$P_{0} = P(|b_{0}| \ge \delta^{2}) < \frac{\epsilon}{3}$$

$$P_{1} = P(b_{1} \ge -\frac{1}{2} J(\theta_{0})) < \frac{\epsilon}{3}$$

$$P_{2} = P(|b_{2}| \ge 2N) < \frac{\epsilon}{3}$$
(3.6.3)

If we let S denote the set of all points  $z^n$  such that,

 $|b_0| < \delta^2$ ,  $b_1 < -J(\theta_0)/2$  and  $|b_2| < 2N$  are all satisfied, then the complementary set S' will consist of all points  $Z^n$ , such that at least one of the inequalities is not satisfied. Therefore, by an elementary law of probability, we have that

 $P(S') \le p_0 + p_1 + p_2 < \varepsilon$ .

Hence,  $P(S) > 1 - \varepsilon$ , and thus  $P(Z^n \epsilon S) > 1 - \varepsilon$  whenever

 $n > n_o(\delta, \varepsilon).$ 

For  $\theta = \theta_0 \pm \delta$ , the right hand side of equation (3.6.2) becomes

$$b_{o} \pm b_{1}\delta + \frac{\lambda}{2}b_{2}\delta^{2}$$

 $b_1 \delta < J(\theta_0)\delta/2$ .

Hence for every point  $Z^{n} \in S$ ,

$$|b_{0} + \frac{\lambda}{2} b_{2} \delta^{2}| < (N + 1) \delta^{2}$$
,

and

Now if  $\delta < J(\theta_0)/2(N+1)$ , the sign of the right hand side of (3.6.2) will be determined by the sign of  $\pm b_1 \delta$ , for  $\theta = \theta_0 \pm \delta$ , hence  $\frac{\partial}{\partial \theta} L_n(\theta) > 0$  for  $\theta = \theta_0 - \delta$  and  $\frac{\partial}{\partial \theta} L_n(\theta) < 0$ for  $\theta = \theta_0 + \delta$ . Furthermore, by RC1, the function  $\frac{\partial}{\partial \theta} L_n(\theta)$ is a.e. $(Z^n)^*$  a continuous function of  $\theta \in \bigoplus$ . Therefore, for sufficiently small  $\delta$  and  $\varepsilon$ , the likelihood equation has a root (in the open interval  $(\theta_0 - \delta, 0_0 + \delta)$ ) with probability larger than  $1-\varepsilon$ , for  $n > n_0(\delta, \varepsilon)$ .

The proof of asymptotic efficiency follows Bar-Shalom (1970).

Let  $\hat{\theta}_n$  be a root of the likelihood equation (3.6.2) in the interval ( $\theta_0 - \delta$ ,  $\theta_0 + \delta$ ). We can arrange equation (3.6.2) as follows

a.e. denotes almost everywhere

$$\hat{\theta}_{n} - \theta_{o} = -\frac{bo}{J(\theta_{o})} / \left[ \frac{b_{1}}{J(\theta_{o})} + \frac{\lambda}{2^{J(\theta_{o})}} (\theta_{n} - \theta_{o}) b_{2} \right].$$
(3.6.4)

Since for  $n > n_0(\delta, \epsilon)$ ,  $|\hat{\theta}_n - \theta_0| < \delta$  with probability exceeding 1- $\epsilon$ . Using (3.6.3) we can show that the denominator of (3.6.4) converges to 1. From RC2, it follows that the numerator of (3.6.4) has mean value zero and variance given by

$$E\left[(b_0/J(\theta_0))^2\right] = \frac{1}{(J(\theta_0))^2} E\left[\left(\frac{1}{n}\sum_{i=1}^n \frac{\partial}{\partial \theta} \log p_i(\theta)\right)^2\right]$$
$$= \left[nJ(\theta_0)\right]^{-2} E\left[\left(\sum_{i=1}^n \frac{\partial}{\partial \theta} \log p_i(\theta)\right)^2\right],$$

which, for  $n > n_0(\delta, \epsilon)$ , is equal to the variance of  $\hat{\theta}_n$ . Hence the total information for the parameter  $\theta$  from n dependent observations is given by

$$J^{n}(\theta_{o}) = (nJ(\theta_{o}))^{2} / \left[ E \begin{bmatrix} \frac{n}{\Sigma} & \frac{\partial}{\partial \Theta} \log P_{i}(\theta) \end{bmatrix}^{2} \\ i=1 & \frac{\partial}{\partial \Theta} \end{bmatrix} \theta = \theta_{o} \quad (3.6.5)$$

It is easily seen from the additivity property of information that the total information is equal to the sum of the information contained in each observation, if and only if condition RC6 is satisfied.

The denominator of (3.6.5) can be written as

$$E\begin{bmatrix}n\\ \Sigma\\ i=1\end{bmatrix}\begin{bmatrix}\frac{1}{2\theta} & \log p_{i}(\theta)\end{bmatrix}^{2} + 2 & \sum_{i  
on using RC6, with  $n > n_{o}(\delta, \epsilon)$ , becomes$$

$$\sum_{i=1}^{n} E \begin{bmatrix} \partial & \log p_{i}(\theta) \end{bmatrix}^{2} = \sum_{i=1}^{n} J_{i}(\theta_{0}) \\ = nJ(\theta_{0}).$$

Thus, for  $n > n_o(\delta, \epsilon)$ , we may rewrite equation (3.6.5) as  $J^n(\theta_0) = n J(\theta_0)$ , and so

$$E\left[(\hat{\theta}_{n} - \theta_{o})^{2}\right] = nJ(\theta_{o})$$
$$= \sum_{i=1}^{n} J(\theta_{o}).$$
$$Q.E.D.$$

#### CHAPTER IV

## LEAST SQUARES

4.1 <u>Introduction</u>. In most scientific investigations it is frequently necessary to determine the values of certain parameters by means of actual measurements. The observations may be made directly on the values of the unknown quantities (dependent variables) or on certain functions of these quantities. In the latter case, the values of the required quantities must be computed from the observed values of the functions. In order to obtain reasonably accurate results, the observations are usually repeated, either in the same way under the same conditions, or in different ways under varying conditions.

Under these circumstances it will be found that different measurements of the same quantity usually give discordant results, the amount of discrepancies varying with the mode of observation. The question, we are faced with now, is how to determine from these discordant results the true values of the required parameters. Because of these discrepancies in the observations, we cannot expect to obtain our parameters with absolute accuracy. All that we can hope for is to obtain those values which are rendered most probable when all the observations are taken into account.

We illustrate the above difficulties with a concrete example. We are required to determine the coefficient of linear expansion of a certain metal rod, based on measurements of the length of the rod at various temperatures.

Temperature in <sup>o</sup> C	Observed Length in m.m.
20	1000.22
40	1000.65
50	1000.90
60	1001.05

Let c denote the required coefficient of linear expansion and  $l_0$  the length of the rod at  $0^{\circ}$ C. The length l of the rod at any other temperature t may be represented by the equation:

 $l = l_0 + ct$  (4.1.1)

Using this model the data can be transformed into the following set of equations:

°	+	20c	=	1000.22	(4.1.2)
lo	+	40c	=	1000.65	(4.1.3)
lo	+	50c	=	1000.90	(4.1.4)
٤	+	60 <b>c</b>	=	1001.05 .	(4.1.5)

We can use any two of the above equations to determine the values of  $l_0$  and c, but these values will depend on the particular pair of equations used, and in general will be different for each pair. For example, we may find from the following pairs of equations, the different values for  $l_0$  and c:

Eq	uatior	<u>1</u> \$	l <sub>0</sub> (m.m.)	c(m.m./ <sup>0</sup> C)
(4.1.2)	and	(4.1.3)	999.790	.02150
(4.1.2)	and	(4.1.4)	999.767	.02267
(4.1.3)	and	(4.1.4)	999.805	.02025
(4.1.4)	and	(4.1.5)	1000.150	.01500
	etc.		etc.	etc.

We are faced with the problem of choosing the best (in the sense of fitting the observations most closely) pair of results. However, in such a situation, we would be disregarding the bulk of the observations thereby losing the majority of the information available. In order to make use of all the data, we choose as our estimates of the parameters those values which make the sum of squared deviations of the observed values from the predicted values a minimum. The name given to this process is the method of Least Squares.

The theory of Least Squares was first discussed by Legendre (1805) and Gauss (1805) who used it as a tool of estimation. Markov and many other authors have since made significant contributions to the general theory.

The underlying assumption in the models of Gauss and Markov is that the error covariance matrix is the identity matrix. Aitken (1934) proposed a model in which the error covariance matrix was a symmetric positive definite matrix, not necessarily the identity. This model was later shown to be equivalent to the model of Gauss and Markov. Rao (1972) developed a unified theory of Least Squares by which any problem involving regression coefficients can be analysed.

The theory of Least Squares, unlike the method of Maximum Likelihood, has no general optimum properties to recommend it in the finite sample case. However, there exists an important class of situations in which it does provide unbiased estimators (which are linear functions of the observations) with minimum variance. We shall also see that, contrary to Kendall and Stuart Vol. II (1960), Least Squares estimators do have asymptotic properties akin to those of Maximum Likelihood estimators.

The solution of the Least Squares Equation is discussed in many standard texts. We shall therefore not discuss its solution explicitly. Instead, we shall discuss properties of Least Squares estimators under different conditions.

# 4.2 Asymptotic Properties of Least Squares

Eicker (1962) considers the linear regression model

$$y = X_n \beta + \varepsilon \tag{4.2.1}$$

where the subscript n denotes the dependence of the matrix X on the sample size. It is assumed that the components of the error vector  $\varepsilon$  are either (a) uncorrelated or (b) independent with  $E[\varepsilon] = 0$  and  $0 < Var(\varepsilon) < \infty$ . The distributions of the components  $\varepsilon_k$  of  $\varepsilon$  are not assumed to be known, nor are

they assumed to be identical. However, it is assumed that the  $\varepsilon_k$  are elements of a certain set F of distribution functions. With these assumptions Eicker (1962) gives conditions on the set F and on the matrix  $X_n$  such that the least squares estimators of the parameters  $\beta_1, \ldots, \beta_q$  are consistent in case (a) or asymptotically normal in case (b).

Cox and Hinkley (1967) consider a similar model. However, they assume that the first column of the matrix  $X_n$ consists of all ones and that the components  $\varepsilon_k$  of the error vector  $\varepsilon$  are independent and identically distributed with known probability density  $f(\varepsilon,\lambda)$ , with zero mean, where  $\lambda$  is an unknown parameter which gives the spread and possibly also the shape of the distribution of f. They then give a condition under which the least squares estimators of  $\beta_1,\beta_2,\ldots,\beta_q$  are asymptotically normal, and then use the fact that the maximum likelihood estimators of  $\beta_1,\ldots,\beta_q$ are asymptotically normal (see Chapter 3), to show the asymptotic efficiency of the Least Squares estimators.

We shall follow Cox and Hinkley (1967) for the discussion of asymptotic efficiency, and Eicker (1962) for the discussion of consistency.

As noted by Cox and Hinkley (1967), the general mean, although included in the linear regression model, is seldom a parameter of primary interest. We therefore can, without loss of generality, assume that the first column of  $X_n$  consist

of all ones and that the other column sums of  $X_n$  are zero. We can also take the parameters  $\beta_2,\ldots,\beta_q$  to be orthogonal to  $\beta_1$ , the general mean.

Since the distributions of the errors are known, we can use the method of Maximum Likelihood to obtain the least squares estimators of the parameters  $\beta_1, \dots, \beta_q$ . The likelihood function L is given by

$$L = \prod_{i=1}^{n} f(\varepsilon_{i}, \lambda) . \qquad (4.2.2)$$

Thus the log-likelihood is given by

$$\begin{split} & \ell = \frac{n}{2} \log f(\varepsilon_{i}, \lambda) = \frac{n}{2} g(\varepsilon_{ij}\lambda), \quad (4.2.3), \text{ say }. \\ & \text{i=1} \\ \text{We shall assume that the range of regularity of } \ell \text{ does not} \\ & \text{depend on the parameters } \beta_{1}, \beta_{2}, \cdots, \beta_{q}. \\ & \text{Let } \mathbf{j}, \mathbf{k} = 1, \dots, q. \quad \text{Then for fixed } \lambda \\ & \frac{\partial \ell}{\partial \beta_{k}} = \sum_{i=1}^{\Sigma} g'(\varepsilon_{i}, \lambda) \frac{\partial \varepsilon}{\partial \beta_{k}} \mathbf{k}, \quad (4.2.4) \\ & \text{and } \frac{\partial^{2} \ell}{\partial \beta_{j} \partial \beta_{k}} = \sum_{i=1}^{n} g'(\varepsilon_{ij}\lambda) \frac{\partial \varepsilon_{j}}{\partial \beta_{j}} \frac{\partial \varepsilon_{k}}{\partial \beta_{j}} \\ & = \sum_{i=1}^{n} g''(\varepsilon_{i}, \lambda) \mathbf{x}_{ij} \mathbf{x}_{ik}, \quad (4.2.5) \\ & \text{since } \varepsilon_{i} = \mathbf{y}_{i} - \frac{q}{k} \mathbf{x}_{ik} \beta_{k} \\ & \text{Thus } \mathbf{E} \left[ -\frac{\partial^{2} \ell}{\partial \beta_{j} \partial \beta_{k}} \right] = \mathbf{x}_{j}^{*} \mathbf{x}_{k} A_{\varepsilon_{i}}, \quad (4.2.6) \\ \end{split}$$

where  $x_{j}$  is the j<sup>th</sup> column of  $X_{n}$  and

$$A_{\varepsilon_{i}} = \int_{|\mathbb{R}^{1}} f(\varepsilon_{i}, \lambda) \frac{\partial}{\partial \varepsilon_{i}^{2}} \log f(\varepsilon_{i}, \lambda) d\varepsilon_{i} .$$

If  $\lambda_1$  is a component of the parameter  $\lambda$  , we have

 $E\left[\frac{\partial^2 k}{\partial \beta_1 \partial \lambda_1}\right] = 0 \quad \cdot$ 

$$\partial_{\beta} \frac{\partial_{j}^{2} \alpha}{j^{2} \lambda_{1}} = - \frac{n}{\Sigma} \frac{\partial}{\partial \lambda_{1}} g'(\varepsilon_{i}, \lambda) \chi_{ij}, \qquad (4.2.7)$$

$$\prod_{i=1}^{n} \chi_{ij} = 0, j = 2, \dots, q \quad \text{we have}$$

and since

Similarly,

we

$$E\left[\frac{\partial^{2} \ell}{\partial \beta_{j} \partial \beta_{1}}\right] = 0, \quad i = 2, \dots, q.$$
  
formation matrix for  $\begin{bmatrix} \beta \\ \beta \end{bmatrix}$  is

Thus the information matrix for  $\begin{vmatrix} p \\ \lambda \end{vmatrix}$  is

 $J\begin{bmatrix} \beta \\ \lambda \end{bmatrix} = \begin{bmatrix} E & 0 \\ 0 & H \end{bmatrix}, \text{ where } E \text{ is the}$ 

information matrix for the parameters  $\beta_1$  and  $\lambda$ , and H

the information matrix for the parameters  $\beta_2, \ldots, \beta_q$ , with

$$H = \begin{bmatrix} X'_{j} X_{k} A_{\varepsilon_{i}} \end{bmatrix}_{h,\ell} \qquad h,\ell = 2,\ldots,q .$$

By the orthogonality condition, the covariance matrix of the Least Squares estimators of  $\beta_2, \ldots, \beta_d$  is

$$\left[ \left[ x_{j} x_{k} \right]_{j,k} \right]^{-1} \quad \text{Var} (\varepsilon)$$

Thus using the fact that the maximum likelihood estimators of  $\beta_2, \ldots, \beta_q$ , have covariance matrix

$$\begin{bmatrix} \begin{pmatrix} x & x \\ y & k \end{pmatrix}_{j,k} \end{bmatrix}^{-1}$$
  
have that the asymptotic efficiency is 
$$\begin{bmatrix} A & Var(\epsilon_i) \\ \epsilon_i \end{bmatrix}$$

1-1

which is independent of the design matrix  $X_n$ .

# Asymptotic Efficiency

are

In the model considered by Eicker (1962) the design matrix  $X_n$  (n x q) is assumed to be of full rank  $q \leq n$ .

The normal equations obtained by minimizing  $\varepsilon'\varepsilon$ 

$$x'_{n} y = x'_{n} x_{n}^{\beta}$$
 (4.2.8)

where  $\hat{\beta} = \hat{\beta}_n$  is the vector of least squares estimators of  $\beta$ . Let  $H_n = X_n^* X_n$ ; by virtue of the rank assumption on  $X_n$ ,  $H_n^{-1}$  exists, and so

$$\hat{\beta} - \beta = H_n^{-1} x'_n \epsilon$$

$$Var(\hat{\beta} - \beta) = H_n^{-1} x'_n V x_n H_n^{-1},$$
(4.2.10)

where  $V = var (\epsilon)$ .

We now state and prove a theorem which gives necessary and sufficient conditions for the consistency of least squares estimators.

<u>Theorem</u>: A necessary and sufficient condition for the least squares estimator  $\stackrel{\sim}{\beta}$  to estimate  $\beta$  consistently on F, is that  $\operatorname{Ch}_{\min n}^{H} \xrightarrow{\rightarrow \infty}$ , Where  $\operatorname{Ch}_{\min n}^{H}$  is the smallest characteristic root of  $H_n$ .

# Proof:

Sufficiency: Since  $E\left[H_n^{-1}X_n' \varepsilon\right] = 0$ , the estimator  $\beta$  is unbiased. Now, the variance of each component of the vector  $H^{-1}X_n' \varepsilon$  tends to zero if and only if  $E\left[\varepsilon'X_nH_n^{-2}X_n' \varepsilon\right]$  tends to zero. Since  $\mathbb{E}\begin{bmatrix} \epsilon_k^2 \\ \epsilon_k \end{bmatrix}$  is assumed constant (i.e., independent of k and n), we have

$$E\left[\epsilon' X_{n} H_{n}^{-2} X_{n}' \epsilon\right] = 0 (tr(X_{n} H_{n}^{-2} X_{n}'))$$
  
= 0 (tr(H\_{n}^{-1}))  
= 0 (1/Ch\_{min} H\_{n}), (4.2.11)

where O(n) = constant x n.

Therefore, since Ch  $_{\text{min}} \ \text{H}_n \ \rightarrow \infty$  , we have

$$\beta \rightarrow \beta$$
 on F.

<u>Necessity</u>: If we choose all  $\varepsilon_k \sim N(0,\sigma^2) \epsilon_F$ 

we see that

$$\overset{\circ}{\beta} - \beta = H_n^{-1} x_n' \varepsilon \sim N_q \left[ 0, \sigma^2 H_n^{-1} \right].$$

Hence the variance of the ith component of  $\beta - \beta$  is given by Var  $(\overset{\sim}{\beta} - \beta)_{i} = \sigma^{2}(H_{n}^{-1})_{ii}$ , which because of consistency must tend to zero for every  $i = 1, \dots, q$ . Therefore  $\operatorname{tr} H_{n}^{-1} = \overset{q}{\underset{i=1}{\Sigma}} \operatorname{Ch}_{i} H_{n}^{-1}$ must tend to zero and so  $\operatorname{Ch}_{i}$  H must tend to infinity

must tend to zero and so Ch H must tend to infinity.

Q.E.D.

## 4.3 Robustness of the Method of Least Squares.

A statistical procedure is said to be robust if it is insensitive to departures from the assumptions which underlie it.

In the case of the general linear model, if the error vector has mean zero and dispersion matrix  $\sigma^2 I_n$  , where  $\sigma^2$ is an unknown constant, then the estimates of the parameters and their variances will remain valid even if the error vector is not normally distributed; i.e., the method of least squares is robust to non-normality. Box and Watson (1962) have studied the robustness of tests on the regression coefficients with respect to the non-normality of the error vector. They mention that for a particular design matrix X(nxp), the variance meansquare ratio  $R_m$  (obtained in comparison of means), and  $R_v$ (obtained in comparison of variances) have distributions which may be approximated by F distributions with modified degrees of freedom. Box and Andersen (1955) showed that the degrees of freedom for the distribution of  $R_m$  are  $\upsilon_1 = \delta_m p$ and  $v_2 = \delta_m (n - p - 1)$  and for  $R_v$  are  $v_1 = \delta_v p$  $v_2 = \delta_v$  (n-p-1), where for mild departures from normality and for moderate numbers of observations,

 $\delta_{\rm m}^{-1} \simeq 1 - (1/n) \beta_1$  and  $\delta_{\rm v}^{-1} \simeq 1 + \frac{1}{2} \beta_2$ , where  $\beta_1$  and  $\beta_2$  are measures of kurtosis, which have values zero when the distribution is normal. It is therefore seen that  $R_{\rm m}$  will be insensitive to non-normality because the corrective

factor is of order  $n^{-1}$ .

When the assumption of uncorrelated homoscedastic errors is not satisfied, the least squares estimate of Y,  $\hat{Y} = (X'X)^{-1}X'Y$ 

is still unbiased; however it no longer has the property of minimum variance. Thus, heteroscedasticity and correlation of the errors still leaves the least squares estimator unbiased. However, the efficiency of such estimators is reduced. It should be pointed out here that the method of least squares is not robust with respect to biased observations.

# 4.4 A Different Computational Technique

There are many sophisticated techniques available for solving the normal equations. These techniques include, singular value decomposition and Householder transformation. We present here an algorithm which finds least squares estimates by minimization of the error sum of squares, rather than by solving the normal equations. The method involves the use of an iterative scheme, and it can also be used to find Ridge estimates.

## Consider the general Gauss-Markov model

#### $Y = X \tilde{Y} + u$

where X is an n x q matrix with  $q \le n$  for which r = r(X)  $\le q$ , E [u] = 0 and Var [u] =  $\sigma^2 I$ .

The least squares estimator  $\widehat{\gamma}$  of  $\gamma$  is obtained by mini-

mizing

$$|y-X\gamma||^2_2 = u'u.$$
 (4.4.1)

Now,

|У -XÝ

$$\Big|_{2}^{2} = y'y - 2y'x'y + y'x'x y.$$
 (4.4.2)

So consider  $||y-X_Y||_2^2 + \alpha ||X||_2^2$ . (4.4.3) It was shown in Chapter I (1.31) that as  $\alpha$  tends to zero the

least squares estimate  $\hat{Y}(\alpha)$  of (4.4.3) converges to the least squares estimate  $\hat{Y}$  of (4.4.2). We shall use this fact and an iterative technique called the Davidon-Fletcher-Powell Method (DFP), (see Walsh (1975)), to solve (4.4.2).

Using (4.4.3), we have,

$$||\mathbf{y} - \mathbf{X}\mathbf{Y}||_{2}^{2} + \alpha ||\mathbf{Y}||_{2}^{2} = \mathbf{Y}' \mathbf{X}'\mathbf{X}\mathbf{Y} + \alpha \mathbf{Y}'\mathbf{Y} - 2\mathbf{Y}'\mathbf{X}'\mathbf{y} + \mathbf{Y}'\mathbf{y}$$
$$= \mathbf{Y}' [\mathbf{X}'\mathbf{X} + \alpha \mathbf{I}]\mathbf{Y} - 2\mathbf{Y}'\mathbf{X}'\mathbf{y} \quad (4.4.4)$$

+уу.

Davidon considers quadratic functions of the form

 $\frac{1}{2} \gamma' G(\alpha) \gamma + \gamma' b + c,$  (4.4.5)

Where  $G(\alpha)$  is symmetric and positive definite;

here  $G(\alpha) = X' X + \alpha I > 0$  for all  $\alpha > 0$ . (4.4.3) is minimized when

 $\frac{\partial}{\partial \gamma} f(\gamma, \alpha) = 2G(\alpha)\alpha - 2X'y = 0, \qquad (4.4.6)$ which gives  $\hat{\gamma}(\alpha) = \left[ G(\alpha) \right]^{-1} X'y$ . However, for n > 3, it is usually difficult to invert  $G(\alpha)$ , and so we use other methods to find  $\hat{\gamma}(\alpha)$  explicitly.

Notation:

Let 
$$g_k = g(\gamma_k) = \frac{\partial}{\partial \gamma} f(\gamma, \alpha) |_{\gamma} =$$

γ<sub>k</sub>

The Method

1/ Set 
$$d_k = -H_k g_k$$
,  
with  $H_1 = I$ . Then  $d_k$  is in the direction of search  
from the current point  $\gamma_k$ .  
2/ Find  $\lambda_k^*$  which minimizes  $f(\gamma_k + \lambda_k d_k)$ .  
3/ Set  $\sigma_k = \lambda_k^* d_k$ .  
4/ Set  $\gamma_{k+1} = \gamma_k + \sigma_k$ .  
5/ Evaluate  $f(\gamma_{k+1})$  and  $g_{k+1}$ , noting that  
 $g_{k+1}$  is orthogonal to  $\sigma_k$ , i.e.  $g'_{k+1} \sigma_k = 0$ .  
6/ Set  $\omega_k = g_{k+1} - g_k$ .  
7/ Set (i)  $A_k = \sigma_k \sigma_k' / \sigma_k' \omega_k$ ,  
(ii)  $B_k = -H_k \omega_k \omega_k' H_k / \omega_k' H_k \omega_k$ ,  
(iii)  $H_{k+1} = H_k + A_k + B_k$ .  
8/ Set  $k = k + 1$ .

9/ Check  $||\sigma_k||$  or  $||d_k||$ , if either satisfies the desired stopping criterion, then stop, else return to Step 1 and repeat the process.

For an n x n system (i.e. n variables), the scheme will converge in at most n steps. Fletcher and Powell recommend that we perform one extra iteration after the apparent minimum is attained. This will help to avoid false minima.

4.4a A Further Application of the Davidon-Fletcher-Powell Method.

We noted in Chapter two that a desirable property of estimators is that of unbiasedness. However, as shown by Hoerl, Kennard and others, in many problems in multiple regression, the Least Squares estimates have a high probability of being unsatisfactory, or even incorrect, especially when the matrix X' X is singular. Hoerl and Kennard (1970) have shown that the addition of a scalar matrix  $kI_p$ ,  $k \ge 0$ , to X' X usually corrects this problem, and for a particular choice of k, we obtain a point estimate of the regression parameter which has a smaller mean square error than than the Least Squares Estimate. This is a particular case of biased estimation and is called Ridge Regression. We note here that in view of (1.31), the Davidon-Fletcher-Powell method can be modified so as to produce Ridge estimates and calculate the mean square error. Example (Goldberger):

Consider the following macro-economic production function

 $y_{t} = \beta_{1} x_{t_{1}} + \beta_{2} x_{t_{2}} + \beta_{3} x_{t_{3}} + \beta_{4} x_{t_{4}} + \varepsilon_{t},$ 

where  $y_t$  is the real gross national product in billions of dollars.  $X_1 = 1$ ,  $X_2$  is labour inputs in millions of man-years,  $X_3$  is real capital in billions of dollars, measured from an arbitrary origin and  $X_4$  is the time in years measured from  $1929 \equiv 1$ , as the base year. The sample consists of 23 annual observations for the United States from 1929 to 1941 and 1946 to 1955.

The observations are as follows.

<sup>x</sup> 2	х <sub>3</sub>	x <sub>4</sub>	ÿ
47	54	1	142
43	59	2	127
39	57	3	113
34	48	4	98
34	36	5	94
36	24	6	102
38	19	7	116
41	18	8	128
42	22	9	140
37	24	10	131
40	23	11	143
42	27	12	157
47	36	13	182
51	9	18	209
53	25	19	214
53	39	20	225
50	51	21	221
52	62	22	243
54	75	23	257
54	94	24	265
55	108	25	276
----	-----	----	-----
52	118	26	271
54	124	27	291

This data was analysed using the Davidon-Fletcher-Powell method, and convergence was achieved after four iterations. The solution is

 $y = -61.6196 + 3.8116 X_2 + .3121 X_3 + 3.8474 X_4$ . The solution given by Goldberger is

 $y = -61.728 + 3.8191 X_2 + .3219 X_3 + 3.7862 X_4$ . Golberger's computation was done using six decimal digit accuracy whereas our computation was done using 15 decimal accuracy and is therefore more accurate. The error sums of squares is 163827.8 . The computations were made on the IBM 360/370 computer of the McGill Computing Center.

Iterate	β <sub>L</sub>	β <sub>2</sub>	β <sub>3</sub>	β <sub>4</sub>	F	<del>∂F</del>   ∂β	σ
1	-50.00000	5.00000	• 0.0	5.00000	133682.0	345052.00	· _
2	-50.01198	4.44023	-0.60235	4.81897	940607.1	33455.75	.000004
3	-50.03816	3.36594	0.26694	4.59982	-300116.2	1145.11	.000065
4	-50.02924	3.48767	0.31627	4.04777	-245180.0	5.91	2.092224
5	-61.61956	3.81155	0.31214	3.84736	163827.8	< 10 <sup>-10</sup>	.7550133

Results for Goldberger's Example

## 4.5 Estimability of Linear Contrast in the General Linear Model (GLM)

The concept of estimability, that is, the existence of unbiased linear estimators for linear contrasts in the general linear model with less than full rank, was introduced by Bose (1944). This concept was developed because the conventional methods, which relied on the inversion of a matrix could not be applied when the matrix is singular. The concept was later discussed by many other authors: Roy and Roy (1959, 1961) who gave a condition for estimability; Searle (1965) who seems to be the first author to use generalized inverses to characterize estimable functions; Miliken (1971) who gives two equivalent conditions for estimability which are specific forms of the condition given by Roy and Roy (1961) and Rao (1972) who developed a unified theory of Least Squares.

In this section we shall discuss some of the many results stated by these and other writers. We shall also prove some theorems on the various characterizations of estimability.

We consider the general linear model

## $y = X\beta + \varepsilon$

where y is an n x l random observation vector, X is an n x p matrix of known constants, of rank q (q < p),  $\beta$  is a p x l vector of unknown parameters defined in R , and  $\varepsilon$  is an unobservable random vector with mean 0 and covariance matrix  $\sigma^2 V$ where  $\sigma^2$  is positive and unknown, and V is a non-negative definite matrix. When V is positive definite, we can reparameterize the model so that the resulting error vector has covariance matrix  $\sigma^2$ I, we shall therefore restrict our discussion to the latter case.

A considerable amount of research has been recently done on the general linear model in which the error vector has a singular covariance matrix. The most important discussions on this model are due to Rao (1971, 1972, 1973a, 1973b, 1973c, 1974, 1975). This model was also studied by several other authors, including Zyskind (1967), Zyskind and Martin (1969), and Alalouf (1975), who,noting that the unified theory developed by Rao (1971, 1972, 1973a) failed to provide an explicit algebraic criterion for a function A'Y to be unbiased for the linear contrast L' $\beta$ , developed a theory similar to Rao's unified theory of Least Squares to handle this problem. We shall present some of the main ideas from Rao (1971, 1972, 1973a, 1973b, 1973c) and Alalouf (1975).

We now state and prove some theorems relevant to the forthcoming discussion.

<u>Definition</u>: The linear contrast  $L'\beta$  is said to be estimable, whenever there exists a linear unbiased estimator  $\Lambda'y$ , such that;

$$E \left[ A' y \right] = L' \beta = A' X \beta \qquad (4.5.1)$$

for all values of  $\beta$ . Then L' = A'X, if and only if L' $\beta$  is estimable.

## Theorem (Roy and Roy (1961)):

The linear contrast  $L'\beta$  is estimable if and only if

$$r(XT) = r(X) - r(L),$$
 (4.5.2)  
 $T \in \mathcal{N}(L)^{*}.$ 

where

Proof: Using (1.26), we have,

$$\mathbf{r}\begin{bmatrix}\mathbf{X}\\\mathbf{L}\end{bmatrix} = \mathbf{r}(\mathbf{L}') + \mathbf{r}\begin{bmatrix}\mathbf{X}(\mathbf{I} - (\mathbf{L}')^{\mathsf{T}}\mathbf{L}')\end{bmatrix}. \qquad (4.5.3)$$

The quantity on the right of equation (4.5.3) is equal to r(X) if and only if

$$r[X(I - (L')]] = r(X) - r(L),$$

but  $I - (L')^{-}L' \in \mathcal{N}(L')$  so  $I - (L')^{-}L'$  is a possible choice of T. We now show that XT is invariant under any choice of T.

Let  $T = I - (L')^{-}L' = UV'$  (a full rank decomposition) where U and V are q x t matrices each of rank

t = r(T) = q - r(L).

Then

r(XT) = r(X U V')  $> r(X U V' V (V'V)^{-1})$  > r(X U)

 $> r(X \cup V') = r(XT)$ .

Hence, r(XT) = r(XU), where the columns of U are a basis for the column space of T. Any other choice of T must admit the columns of U as a basis. Q.E.D. \* $\mathcal{N}(\cdot)$  denotes the null space. Theorem I (Miliken (1971)).

The linear contrast L'  $\beta$  is estimable, if and only if  $r \left[ X(I-(L')^{+}L') \right] = r(X) - r(L)$ . (4.5.4) This theorem is equivalent to the theorem of Roy and Roy (1961) since I - (L')^{+} L' $\epsilon \mathcal{N}(L')$ . <u>Theorem II (Miliken (1971))</u>. The linear contrast L'  $\beta$  is estimable if and only if  $tr \left[ X(I-(L')^{+}L') \{ X(I-(L')^{+}L') \}^{+} \right] = q-t$ . (4.5.5) <u>Proof</u>: Note that the matrix in (4.5.4) is idempotent, hence

$$tr \left[ X(I-(L')^{+}L') \{ X(I-(L')^{+}L') \}^{+} \right] = r \left[ X(I-(L')^{+}L') \{ X(I-(L')^{+}L') \}^{+} \right]$$
$$= r \left[ X(I-(L')^{+}L') \right]$$
$$= r(X) - r(L)$$
$$= q - t .$$
Q.E.D.

Theorem (Searle (1965)).

The linear contrast  $L' \beta$  is estimable if and only if

L' X - X = L' (4.5.6)

for some and hence for every  $X^- = g_1(X)$ .

<u>Proof</u>; If L' = L'X  $\overline{X}$  then choosing A' = L'X, we have L' = A'X, thus L'  $\beta$  is estimable. Conversely, L' X<sup>-</sup> X = A' X X<sup>-</sup> X, and putting L' = A'X

= A'X= L'.

Q.E.D.

Theorem: Let M be a symmetric matrix such that

$$\mathcal{G}(X'MV) \subset \mathcal{G}(X'MX)$$
, in which case

 $(y-X\beta)'M(y-X\beta)$  , as a function of  $\beta$  , has stationary values.

Let  $\hat{\beta}$  be a stationary point.

If  $l'\hat{\beta}$  is the BLUE of  $l'\beta$  for every  $l \in \mathcal{C}(X)$  then it is necessary that

$$r(X'MX) = r(X)$$
 (4.5.10)

and M is of the form

$$(V + XUX) + K$$
 (4.5.11)

for any symmetric choice of g- inverse, where U and K are symmetric matrices such that

$$G_{\rm Q}(V;X) = G_{\rm Q}(V + XUX'),$$
 (4.5.12)

$$VKX = 0, X'KX = 0.$$
 (4.5.13)

Conversely if M is of the form (4.5.11) with 4.5.12) and (4.5.13), then r(X'MX) = r(X) and  $l'\hat{\beta}$  is the BLUE of  $l'\beta$ for every  $l\in Q(X')$ . <u>Proof</u>: Following Rao (1972), we have, on equating the

derivative of  $(y - X\beta)'(y - X\beta)$  to zero

$$X' MX \beta = X' My$$
 (4.5.14)

This system is consistent since  $\int_{Q} (X'MV) C \int_{Q} (X MX)$ 

and  $y \in G(V; X)$ , w.p. 1, hence

$$B = (X'MX) X'My$$
 (4.5.15)

is a stationary point. Let  $\ell$  = X'L. Then  $\ell'\hat\beta$  is the BLUE of  $\ell'\beta$  , hence by Searle's theorem

$$L'X(X'MX) X'MX = L'X.$$

Since L is arbitrary, we have

$$X(X'MX)^{T}X'MX = X, \text{ hence}$$
(4.5.16)  
$$r(X) = r \left[ X(X'MX)^{T} X'MX \right] \leq r(X'MX) \leq r(X).$$

i.e., r(X'MX) = r(X), which proves (4.5.10). If  $l'\hat{\beta}$  is the BLUE of  $l'\beta$  for every  $l \in Q(X')$ , then applying the lemma on pp. 317 of Rao (1973a), we have

$$L'X(X'MX) X'MVZ = 0$$
 (4.5.17)

for any L, where Z is a matrix of maximum rank such that  $Z \in \mathcal{N}(X')$ . Therefore (4.5.17) implies that

$$X(X'MX) - X'MVZ = 0$$

if and only if X' MVZ = 0, this implies that

$$VMX = XP \qquad (4.5.18)$$

for some P. Now, there exists a symmetric matrix U, such that

$$X'M(V + XUX')MX = X'MX$$
 (4.5.19)

One choice of U is (X'MX) (X'MX - X'MVMX)(X'MX), since X'MXUX'MX

$$= X'MX(X'MX) X'MX(X'MX) X'MX - X'MX(X'MX) X'MVMX(X'MX) X'MX$$

= X'MX - X'MX(X'MX) X'MXP(X'MX) X'MX

= X'MX - X'MXP(X'MX) X'MX



= x'Mx - x'MVMX(x'MX) x'MX

= X'MX - X'MVMXby (4.5.16).

Multiplying both sides of (4.5.19) by X(X'MX) and using (4.5.8) and (4.5.18), we obtain

$$(V + XUX')MX = X.$$
 (4.5.20)

If p'(V + X U X') = 0, then from (4.5.20) p'X = 0, and hence p'V = 0 and vice versa, proving (4.5.12). Choosing a g<sub>12</sub> (V + X U X') and a symmetric matrix K, let

$$M = (V + XUX') + K, \qquad (4.5.21)$$

substituting (4.5.21) in (4.5.20), we obtain (4.5.13). The converse is easily verified by using lemma 2.

> Let  $G = X'X + \Sigma = X'X + \Gamma\Gamma'$  where  $G = g_1(G)$  $S = X'G^T$  $H = H_0 + H_1$

where

 $H_0 = (I - SS^+)X'G^-, H_1 = (S^+)'F'G^-.$ Alalouf (1975) states and proves a lemma which gives the properties of the matrices G, S,  $H_0$  and  $H_1$ . We shall use these results in the sequel and refer the reader to Alalouf (1975) for the proofs. He then uses these results to obtain the following decomposition of the vector space into four virtually disjoint subspaces in terms of four projectors;

$$\begin{array}{cccc} \underline{Projector} & \underline{Subspace} & \underline{Dimension} \\ \hline XH_1 & & & & & \\ & & & \\ XH_0 & T_1 \subset & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ &$$

He then gives the following theorem which justifies the above decomposition.

<u>Theorem</u>. Let G, H,  $H_0$ ,  $H_1$  be as defined above. Then

(i) 
$$H_0 Y = H_0 X\beta$$
 w.p.1. (4.5.22)  
(ii)  $Cov \left[ XH_1 Y, (GG^- - XH)Y \right] = 0, Cov \left[ XHY, (GG^- - XH)Y \right] = 0.$   
(iii)  $(I - GG^-) Y = 0$  with probability one (w.p.1).

Using the decomposition and the previous theorem we obtain the following decomposition of y;

$$y = XH_0 y + XH_1 y + (GG - XH)y + (I - GG)y$$
 (4.5.23)

and note that the fourth term is identically zero and so provides us with no information, while the first and second terms provide information about  $\beta$ . The third term has mean zero and so gives no information about  $\beta$ .

The following lemma gives an expression for the permissible values of  $\beta$ .

<u>Lemma</u>. The vector  $\beta$  satisfies

$$\beta = H_0 y + (I - H_0 X) \alpha, \qquad (4.5.24)$$

where  $\alpha$  is an arbitrary vector and  $H_0^y$  is fixed w.p.1 <u>Proof.</u> Using parts (iv) and (vi) of Lemma 5.1 (Alalouf) we have that (4.5.24) is the general solution of (4.5.22). <u>Theorem.</u> A'y is an unbiased estimator of L' $\beta$ if and only if

$$A'X[I - H_0X] = L'[(I - H_0X]],$$
 (4.5.25)

and

$$A'XH_0 \dot{y} = L'H_0 \dot{y}$$
 (4.5.26)

Proof. If A'y is unbiased for L'
$$\beta$$
, then for all  $\beta$   
satisfying (4.5.24), we have, A'X $\beta$  = L' $\beta$ .  
Thus A'XH<sub>0</sub>y + A'X [I - H<sub>0</sub>X]  $\alpha$  = L'H<sub>0</sub>y + L'[I - H<sub>0</sub>X]  $\alpha$   
must hold for every  $\alpha$ . Hence (4.5.25) and (4.5.26).  
Conversely, if A' and L' satisfy (4.5.25) and(4.5.26),  
then

$$E [A'y] = A'X\beta = A'XH_0y - A'X[I - H_0X]\alpha$$
$$= L'H_0y - L'(I - H_0X)\alpha$$
$$= L'\beta.$$

Q.E.D.

## BIBLIOGRAPHY

- Aitchison, J., and Silvey, S. D. (1958), "Maximum Likelihood Estimation of Parameters Subject to Restrictions", <u>Ann Math. Statist</u>. 29, 813-828.
- Alalouf, I. S (1975), "Estimability and Testability in Linear Models", Ph.D. Thesis, McGill University.
- Anderson, T. W. (1958), "An Introduction to Multivariate Statistical Analysis", New York: John Wiley & Sons.
- Anderson, T. W. (1961), "Least Squares and Best Unbiased Estimates", <u>Ann. Math. Statist</u>. 33, 266-273.
- Bahadur, R. R. (1958), "Examples of Inconsistent Maximum Likelihood Estimates", Sankhyā, 20, 207-210.
- Bar-Shalom, Y. (1970), "On the Asymptotic Properties of the Maximum Likelihood Estimate Obtained from Dependent Observations",

J.R. Statist. Soc. B., 33, 72-77.

Barnett, V. D. (1966), "Evaluation of the Maximum Likelihood Estimator where the Likelihood Equation has Multiple Roots", <u>Biometrika</u> 53, 1 and 2, 151-165.

Bartlett, D. P. (1915), "The Method of Least Squares", <u>Rumford Press</u>. Berkson, J. (1956), "Estimation by Least Squares and by Maximum

Likelihood"., <u>Proceedings of the Third Berkeley Symposium on</u> <u>Mathematical Statistics and Probability, Berkeley and Los</u> Angeles, University of California Press, 1956, pp. 1-11. Bhat, B. R. (1974), "On the Method of Maximum Likelihood for De-

pendent observations", <u>J.Roy. Statist. Soc. B</u>, 36, 48-53. Box, G. E. P., and Draper, N. R. (1975), "Robustness to nonnormality

of regression tests", Biometrika 49, 1 and 2, 93-106.

Chambers, J. M. (1975), "Fitting non-linear models: numerical techniques", <u>Biometrika</u> 60, 1, 1-10.

Comstock, G. C. (1889), "Method of Least Squares", <u>Ginn and Company</u>. Conte, S. D., and deBoor, C. (1972), "Elements of numerical analysis: An algorithmic approach", 2nd ed., <u>McGraw-Hill</u>.

- Cox, D. R., and Hinkley, D. V. (1966), "A note on the efficiency of Least Squares Estimates", <u>Ann. Math Statist</u>. 37, No. 2, 284-289.
- Cramér, H. (1946), "Mathematical Statistics", Princeton University Press.

Daniels, H. E. (1961), "The asymptotic efficiency of a Maximum Likelihood Estimator", <u>Proceedings of the Fourth Berkely Symposium</u> <u>on Mathematical Statistics and Probability, Berkeley and Los</u> <u>Angeles</u>, University of California Press, Vol. 1, p. 151-163.

Draper, N., and Smith, H. (1966), "<u>Applied Regression Analysis</u>", New York: Wiley & Sons.

Edwards, A. W. F. (1972), "Likelihood", Cambridge University Press.

Eicker, F. (1962), "Asymptotic Normality and Consistency of the Least Squares Estimators for families of Linear Regressions", <u>Ann</u>. <u>Math. Statist. 34, 447-456</u>.

- Fisher, R. A. (1950), "Contributions to Mathematical Statistics", New York: Wiley & Sons.
- Golub, G. H., and Styan, G. P. H. (1973), "Numerical Computations for Univariate Linear Models". J. Statist. Comput. Simul. Vol. 2 253-274.
- Graybill, F. A. (1961), "An Introduction to Statistical Models, Vol. 1, McGraw-Hill.
- Graybill, F. A., and Marsaglia, G. (1957), "Idempotent matrices and Quadratic Forms in the General Linear Hypothesis", <u>Ann. Math</u>. <u>Statist.</u> 28, 678-686.
- Hoerl, A. E., and Kennard, R. W. (1970a), "Ridge Regression: Biased Estimation for Non-Orthogonal Problems", <u>Technometrics</u>, Vol. 12, No. 1, 55-68.
- Hoerl, A. E., and Kennard, R. W. (1970b), "Ridge Regression: Applications to Nonorthogonal Problems" <u>Technometrics</u>, Vol. 12, No. 2 69-82.
- Hoerl, A. E., Kennard, R. W., and Baldwin, K. F. (1974), "Ridge Regression: Some Simulations", <u>Communications in Statistics</u>, 4(2), 105-123.
- Hoadley, B. (1971), "Asymptotic Properties of Maximum Likelihood for the Independent not Identically distributed Case". Ann. Math. Statist. 42, 6, 1977-199.
- Hogg, R. V., and Craig, A. T. (1961), "Some Results on Unbiased Estimation", <u>Sankhya</u> A, 24, 333-338.

Kale, B. K. (1961), "On the Solution of the Likelihood Equation by Iteration Processes", <u>Biometrika</u> 48, 452-456.

- Kale, B. K. (1962), "On the Solution of Likelihood Equations by Iteration Processes: The Multiparameter Case", <u>Biometrika</u> 49, 3 and 4, 479-486.
- Kendall, M. G., and Stuart, A. (1967), "The Advanced Theory of Statistics", Vol. 2, Hafner.
- Lewis, T. O., and Odell, P. L. (1966), "A Generalization of the Gauss-Markov Theorem", J.A. S. A. 66, 1063-1066.
- Marsaglia, G., and Styan, G. P. H. (1974), "Equalities and Inequalities for Ranks of Matrices", <u>Linear and Multilinear</u> Algebra, Vol. 2
- Milliken, G. A. (1971), "New Criteria for Estimability for Linear Models", Ann. Math. Statist. 42, 5, 1588-1594.
- Norton, H. W. (1956), "One Likelihood Adjustment may be Inadequate", Biometrics 12, 79-81.
- Prasad, M. S., and Rao, B. L. S. P. (1976), "Maximum Likelihood Estimation for dependent Random Variables", <u>J.I.S.A</u>. Vol. 14, 79-97.
- Rao, C. R. (1945), "Generalizations of Markov's Theorem and Tests of Linear hypotheses", Sankhya 7, 1, 9-19.

Rao, C. R. (1946), "On the linear combination of observations and the general theory of Least Squares", <u>Sankhya</u> 7, 3, 237-256.
Rao, C. R. (1950), "A theorem on Least Squares", <u>Sankhya</u> 11, 1, 9-12.

- Rao, C. R. (1958), "Maximum Likelihood Estimation for the multinomial distribution with infinite numbers in cells", <u>Sankhya</u> 20, 3 and 4, 211-218.
- Rao, C. R. (1972), "Unified theory of Least Squares", <u>Communications</u> in Statistics 1(1), 1-8.
- Rao, C. R. (1973a), "Linear Statistical Inference" 3rd edition. New York: Wiley and Sons.
- Rao, C. R. (1973b), "Representations of best linear unbiased estimators in the Gauss-Markov model with a singular dispersion matrix", J. Multivariate Analysis, 3, 276-292.
- Rao, C. R. (1973c), "Projectors, Generalized Inverses and BLUE'S", J. Roy. Statist. Soc. B, 36, 442-448.
- Rao, C. R., and Mitra, S. K. (1971), "Generalized inverses of matrices and its applications", New York: <u>Wiley & Sons</u>,
- Searle, S. R. (1964), "Additional results concerning estimable functions and generalized inverse matrices", <u>J. Roy. Statist. Soc. B</u>, 27, 486-490.
- Seber, G. A. F. (1977), "Linear Regression Analysis", New York, Wiley & Sons.
- Silvey, S. D. (1961), "A note on Maximum Likelihood in the case of dependent random variables", <u>J. Roy. Statist. Soc. B</u>, 23, 444-452.
- Wald, A. (1949), "Note on the consistency of the Maximum Likelihood Estimate", <u>Ann. Math. Statist</u>. 20, 595-601.

Walsh, G. N. W. (1975), "Methods of Optimization", London: Wiley & Sons.

Wolfowitz, J. (1948), "On Wald's proof of the consistency of the Maximum Likelihood", <u>Ann. Math. Statist</u>. 20, 601-602. Zyskind, G., and Martin, F. B. (1969), "On best linear estimation

and a general Gauss-Markov theorem in linear models with arbitrary non-negative covariance structure", <u>SIAM. J</u>. <u>Appl. Math. Vol. 17, 6, 1190-1202.</u>