

Mimeo Series

May 96

No. 2161T

Statistical Estimation Under  
Length Biased Distributions

by

Stoffel R. Moeng

Dept. of Biostatistics

Univ. of North Carolina

---

Name

Date

---

**STATISTICAL ESTIMATION UNDER  
LENGTH BIASED DISTRIBUTIONS**

by

**Stoffel R. Moeng**

Department of Biostatistics

University of North Carolina

Institute of Statistics

Mimeo Series No. 2161T

May 1996

# STATISTICAL ESTIMATION UNDER LENGTH BIASED DISTRIBUTIONS

by

Stoffel R. Moeng

A dissertation submitted to the faculty of the University of North Carolina at Chapel Hill in partial fulfillment of the requirements for the degree of Doctor of Philosophy in the Department of Biostatistics

Chapel Hill

1996

STOFFEL R. MOENG. Statistical Estimation Under Length-Biased Distributions. (Under the joint direction of Drs. P.K. Sen and C.M. Suchindran.)

### ABSTRACT

In simple random sampling, there is an equal chance of selection for each unit in the population. However, there are some practical situations where this might not be ideal. We consider the case of length-biased sampling, whereby the chance of selection of a unit is proportional to the length of the unit, by which the length-biased density is defined. Specifically we consider the case of a continuous random variable  $Y$  with pdf  $f(y;\theta)$ , where,  $f(y;\theta)$  is of an exponential family of distributions and  $\theta$  is a  $k$ -parameter vector to be estimated. We show the bias of the MLE of  $\theta$  to be of order  $n^{-1}$  and use the jackknife estimation technique to eliminate the leading term in the bias of the MLE of  $\theta$ . It is shown that the Jackknife estimator has the same normal distribution as the MLE in large samples.

The estimation study is extended to the regression problem where the mean of the sufficient statistics of the exponential family of distributions depends on a set of covariates. Both the bias and asymptotic distribution of the MLE and the Jackknife are provided. The Fisher information from the length-biased distribution and the original distribution are compared through the A-optimality and D-optimality which are functions of the eigenvalues of the Fisher information matrices.

Data from the Demographic and Health Surveys is analyzed by maximum likelihood methods and Jackknife methods. The response variables is the first order birth interval and the potential covariates are the age at marriage and duration of marriage for the mothers. Simulation studies from the lognormal and gamma distributions indicate a reduction in the bias of the MLE's for the scale parameter of the lognormal distribution. Both the MLE and the Jackknife estimates are approximately normal. The RMSE for the MLE and the Jackknife estimates are given.

## ACKNOWLEDGEMENTS

First of all, I would like to thank professor P.K. Sen for his patient guidance and all his advice in the preparation of this dissertation. I would also like to thank Professor C.M. Suchindran for his joint direction and his many helpful encouragement, comments and suggestions on how best to shape this material. I am grateful to my other dissertation committee members, professor M. Symons, Professor G. Koch and Professor R. Bilsborrow for their participation in my research adventure. I owe a particularly deep debt of gratitude to Professor Bilsborrow for his efforts "above and beyond the call of duty" during my studies in helping me realize this goal.

I am grateful for the support of my friends, the Hogan family who shared with me the love of their home and with whom I have shared Thanksgivings; my fellow graduate students for their friendship and encouragement during my education. of course, none of this effort would have been possible without the love and support of my family; my parents, aunts, uncles, and sisters. My special thanks to my aunt Dorah who died just after I had started my studies at UNC and well after she had left her marks upon my life.

I am deeply thankful to the entire faculty of the Department of Biostatistics at UNC for their constant support and motivation during my studies. Finally I would like to acknowledge the university of North Carolina for the opportunity to attend and learn a valuable education.

## TABLE OF CONTENTS

Chapter	Page
I. INTRODUCTION AND LITERATURE REVIEW.....	1
1.1. Introduction .....	1
1.2. Design Issues .....	4
1.3. Literature review .....	5
1.3.1 Origins and Applications .....	5
1.3.2 Relationship between Unweighted and Weighted Distributions.....	8
1.3.3 Comparison of Unweighted and Weighted experiments.....	11
1.3.4 Estimation of Parameters .....	13
1.3.5 Density Estimation.....	19
1.3.6 Discussion and Summary of Literature review .....	20
1.4 Proposed Research.....	22
II. PARAMETER ESTIMATION: EXPONENTIAL FAMILY .....	26
2.1. Introduction .....	26
2.2. Maximum Likelihood Estimators.....	26
2.2.1. Asymptotic Distribution and Bias of the MLE.....	32
2.3. Jackknife Estimators of the MLE .....	41
2.3.1. Asymptotic Distribution and Bias of the Jackknife Estimator.....	47
III. PARAMETRIC REGRESSION MODEL: EXPONENTIAL FAMILY .....	51
3.1. Introduction .....	51
3.2. Estimation of Parameters of the Regression Model .....	56
3.2.1. Bias of the Regression Parameter Estimators.....	66
3.3. Jackknife Estimator of the Regression Model Parameters .....	71

IV. LENGTH-BIASED DISTRIBUTION IN EXPONENTIAL FAMILY .....	78
4.1. Introduction .....	78
4.2. Maximum Likelihood Estimators of the Parameters .....	79
4.3. Asymptotic Distribution and Bias of the MLE.....	80
4.4. Estimating Regression Parameters from the Length-Biased Distribution.....	86
V. EXAMPLE AND SIMULATION STUDY.....	91
5.1. Introduction .....	91
5.2. Results and Discussion .....	94
5.3. Simulation Study.....	103
VI. DISCUSSION AND SUGGESTIONS FOR FUTURE RESEARCH.....	111
BIBLIOGRAPHY .....	113

## LIST OF TABLES

Table5.1:	Mean, Bias and RMSE of the MLE and Jackknife Estimates for M = 1000 samples of size $n$ from a lognormal Distribution Model. 1990 Sudan DHS birth interval data .....	94
Table5.2:	Mean, Bias and RMSE of the MLE and Jackknife Estimates from the Lognormal Regression Model. 1990 Sudan DHS birth interval data. M = 1000 samples of size $n = 50, 100, 200$ . .....	99
Table5.3:	Mean, Bias and RMSE of the MLE and Jackknife Estimates for M = 1000 samples of size $n$ from a Lognormal Distribution. Simulations with $\mu = 2.0$ and $\sigma^2 = 1.5$ .....	105
Table5.4:	Mean, Bias and RMSE of the MLE and Jackknife Estimates for M = 1000 samples of size $n$ from a Gamma Distribution. Simulations with $\mu = 1.5$ and $\sigma^2 = 1.5$ .....	105
Table5.5:	Mean, Bias and RMSE of the MLE and Jackknife Estimates of the Parameters of the Lognormal Regression Model. M = 1000 simulations of sample size $n = 50, 100, 200$ . .....	107

## LIST OF FIGURES

Figure 5.1:	Maximum Likelihood and Jackknifed Estimates of the Scale Parameter of the Lognormal Distribution for the Sudan First Order Birth Interval Data. Sudan DHS (1990) .....	97
Figure 5.2:	Maximum Likelihood and Jackknifed Estimates of the Intercept Parameter of the Lognormal Regression Model for the Sudan First Order Birth Interval Data. Sudan DHS (1990) .....	101
Figure 5.3:	Maximum Likelihood and Jackknifed Estimates of the Marital Duration Parameter of the Lognormal Regression Model for the Sudan First Order Birth Interval Data. Sudan DHS (1990).....	101
Figure 5.4:	Maximum Likelihood and Jackknifed Estimates of the Age Parameter of the Lognormal Regression Model for the Sudan First Order Birth Interval Data. Sudan DHS (1990) .....	102
Figure 5.5:	Maximum Likelihood and Jackknifed Estimates of the Scale Parameter of the Lognormal Regression Model for the Sudan First Order Birth Interval Data. Sudan DHS (1990) .....	102
Figure 5.6:	Maximum Likelihood and Jackknifed Estimates of the Intercept Parameter of the Lognormal Regression. The mean estimate from $M = 1000$ samples of size $n = 50$ .....	109
Figure 5.7:	Maximum Likelihood and Jackknifed Estimates of the Beta_1 Parameter of the Lognormal Regression. The mean estimate from $M = 1000$ samples of size $n = 50$ .....	109
Figure 5.8:	Maximum Likelihood and Jackknifed Estimates of the Beta_2 Parameter of the Lognormal Regression. The mean estimate from $M = 1000$ samples of size $n = 50$ .....	110
Figure 5.9:	Maximum Likelihood and Jackknifed Estimates of the Scale Parameter of the Lognormal Regression. The mean estimate from $M = 1000$ samples of size $n = 50$ .....	110



# Chapter I.

## Introduction and Literature Review.

### 1.1. Introduction.

Statistical inference methods for analysing data are aimed at estimating the parameters of a postulated underlying probability model generating the data. It is usually assumed that the investigator has access to the data following the distribution from which inference is to be drawn. This assumption is not always justified. In statistical analysis, one often has observations from a weighted distribution rather than from the original (unweighted) distribution, and from these observations one needs to make statistical inferences about the parameters of the original distribution.

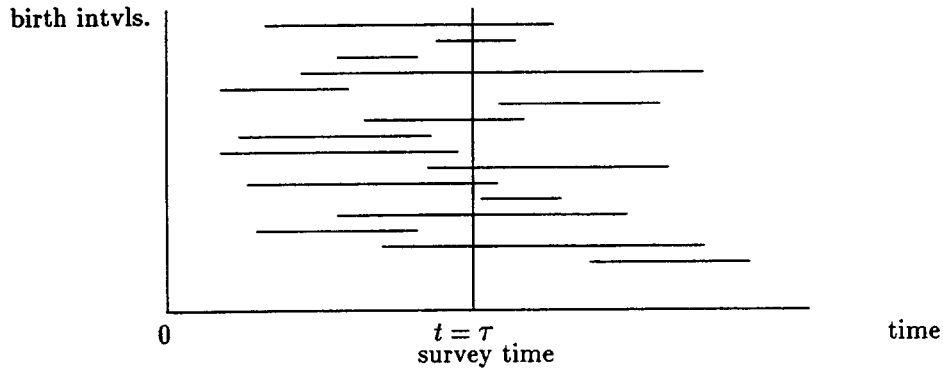
This research is motivated by the problem of statistical estimation of the population parameters from the data generated by a sampling design which gives rise to data from a length-biased distribution, which is a form of weighted distribution. Weighted distributions arise when realizations  $y_i$  of a random variable  $Y$  having a probability density function  $f(y)$  are observed and recorded according to a probability density  $g(y)$  which is a weighted version of  $f(y)$ . Rao (1965) identified various sampling schemes that could be modeled by weighted distributions.

Length-biased sampling has implications for survey data collected at an

arbitrary point in time. Suppose that a component operating in a system is replaced upon failure by another component having the same life distribution, so that the sequence of component life lengths  $\{Y_i\}$  form a renewal process. In order to determine the mean life of the component, components in operation in a particular system are sampled at a fixed time on a particular day and the total lives of these components are measured. Because the components selected are in operation on that particular date, they constitute a nonrandom sample of all components, having longer mean lifetimes than the class of components from which they were drawn. The probability that the sampled component has life length  $y$  is a length-biased version of the probability distribution of the lengths of all the population of life lengths of the components.

The length-biased design is illustrated in Figure 1 for a hypothetical survey data on birth intervals of a given order  $k$ , for a cohort of women married at the same time. Each line represents the duration of time between any two live births for each woman, such that the left endpoint is the time the woman experienced the event of birth, and the next birth occurred at the right endpoint of the line. Suppose at an arbitrary point in time (say  $t = \tau$ ), a survey is undertaken to estimate the mean birth interval for all birth intervals of the given order in the population from which these birth intervals are drawn. The investigator ascertains and records the time of the last birth before the survey and the time the next birth occurs for the births covering the  $k^{\text{th}}$  birth interval. Since longer birth intervals are more likely than shorter birth intervals to overlap the survey time and hence be sampled, the sample of birth intervals thus obtained comprises of a length-biased sample of all  $k^{\text{th}}$  order of birth intervals.

Figure 1. Hypothetical data layout of the  $k^{\text{th}}$  order birth intervals.



Length-biased sampling is often encountered in event history data collection. In renewal event, such as open birth intervals in demographic studies, long intervals are usually over-represented, resulting in biased estimates of the mean birth interval of the population. Survey data taken at an arbitrary point in time to measure the time between the occurrence of an event and the next occurrence of the event tend to draw a length-biased data from the population of all such duration times. Examples of such survey questions include questions such as “when did you last see a physician?” and “when did you last give live birth?”, and the duration of time since the last occurrence and the next occurrence is recorded form a length-biased data (Morrison, 1973).

The length-biased distribution is frequently appropriate for certain natural sampling plans in survival analysis, biometry and reliability (see Zelen, 1974; Gupta and Keating, 1986; Gupta and Kirmani, 1980). Section 1.2 reviews the relevant literature in the area of length-biased distributions.

Important factors in data analysis include the choice of an appropriate model for the data, which has implications for the estimation techniques, that is whether

parametric methods, for example, maximum likelihoods, the method of moments or other nonparametric methods, would be appropriate. Failure to incorporate the selection mechanism in analysing such data leads to biased results.

It is also important to consider the effects of possible covariates on the parameters to be estimated. Information on covariates affecting the response would improve the precision of parameter estimators and provide greater validity to statistical inference. Such analysis is usually carried out by regression analysis methods.

The problems of interest to us are: (i) To obtain unbiased and efficient estimators of the parameters of the original distribution, given that the data are from a length-biased distribution, (ii) To study the asymptotic distribution of such estimators and (iii) To examine the information about the the parameters lost or gained by not having access to the sample following the original distribution.

The objectives of the study may focus on recovering the characteristics of the original unweighted distribution based on a lenth-biased sample, and assessing the nature of distortion in determining these characteristics in case the inherent biased sampling is ignored in modelling the outcome variable. One may also be interested in the estimates of the length-biased distribution in itself.

## 1.2 Design Issues

The length biased or weighted sampling procedure occurs when a sample of units is drawn from a population such that the probability of a unit being selected in the sample is proportional to its length. To motivate the estimation problem, we state the problem as follows:

Let  $\{Y_i; i = 1, \dots, n\}$  be a sequence of independent random variables with

common pdf  $f(y;\theta)$  for  $\theta \in \Omega$ . Instead of obtaining sample observations from the original density  $f$ , the realizations made on the random variable  $Y$  are obtained from the length-biased density

$$g(y;\theta) = \frac{w(y)f(y;\theta)}{\mu(\theta)}, \quad w(y) > 0. \quad (1.2.1)$$

where  $\mu(\theta) = E[w(Y)] = \int w(y) f(y;\theta) dy$ , for some nonnegative weight function  $w(y)$ .

Among the weight functions commonly used are  $y$ ,  $y^\gamma$  ( $\gamma > 0$ ) and  $1-(1-\beta)^y$  for  $0 < \beta < 1$ . The pdf  $g(y;\theta)$  is sometimes referred to as the weighted distribution. As a special case, this research will restrict itself to length-biased distributions where the weight function  $w(y)=y$ . The random sample from length-biased distribution is called a length-biased sample (Cox, 1969 ; Patil and Rao, 1977).

The problem then is to estimate the parameter  $\theta$  or functions thereof on the basis of length-biased sample from density  $g$ .

## 1.3. Literature Review and Related Research.

### 1.3.1. Origins and Applications.

The concept of weighted distributions was first formulated in a unified manner by Rao (1965) as an extension of the basic idea originally developed by Fisher (1934). Rao distinguished three major reasons, or situations, in which a sampling scheme would yield a sample that may be a biased representation of the population, and introduced the concept of weighted distributions as an adjustment

to the original unweighted distribution. Applications of the weighted distributions have since received attention in the literature. Briefly, the three situations identified by Rao are as follows.

i) Non observability / Visibility bias

This occurs frequently in ecological, aerial surveys when a group gets recorded only when at least one member in the group is sighted and each individual has an independent chance  $\beta$  of being sighted. If  $X$  is the number of individuals in a group with pdf  $f(x)$ , the probability that an observed group has  $x$  individuals is a weighted distribution with weight  $w(x) = 1 - (1 - \beta)^x$ . In this case  $E[w(X)]$  is the probability of observing a group. Rao shows that as  $\beta \rightarrow 0$ , the corresponding distribution has a length biased pdf. Rao (1965) discussed the application in the study of the distribution of albino children (individuals with rare case). A convenient sampling method is to first discover an albino child and through it obtain the albino count  $X^w$  of the family to which it belongs. The pdf of  $X^w$  is a length-biased version of the original r.v.  $X$ , the number of albino children in families with proneness to produce such children, similarly, in some instances we may be able to observe patients who have a certain disease that is manifested only when some symptoms are visible and we are interested in the distribution of the entire population.

ii) Partial destruction of observations/ partial ascertainment.

This situation arises when a sample is obtained from a pdf  $f(x)$ , but an observation  $X = x$  is reduced to  $y$ , by a destruction process, with conditional pdf  $\zeta(y/x)$ . Suppose that an observation is recorded only when the original value is unchanged, and the probability that the observation  $X = x$  is undamaged is  $\zeta(x/x)$ . Then the pdf of an observation  $X^w$  is a weighted distribution with a

weight function  $w(x) = \zeta(x/x)$ . In observing human disease, the probability of detection can depend on the length of time the individual has had the disease. Zelen (1974) discussed how health-screening programs tend to sample a length biased distribution of all individuals with the disease, in which the weight is the length of time the subject is in the pre-clinical state and the disease is undetectable. Specifically Zelen argued that women with breast cancer detected in a cancer screening program are a length biased sample of all women with breast cancer. The weight is the length of time that cancer is in the detectable state. Since patients with a short pre-clinical stage may never be enrolled in the screening program, resulting in failure to report some cases, the sample tends to be a biased sample. Patients with slow-growing, less aggressive forms of the cancer will be over represented as their pre-clinical stage is longer.

iii) Sampling with unequal chances of selection.

This design is common in sampling surveys where the investigator intentionally samples units according to a probability proportional to the size (pps) of some measure. For example, let  $(X, Y)$  be a pair of nonnegative random variables with joint pdf  $f(x, y)$  and  $X$  be the main r.v. under study. An observation is made on  $X$  using the conditional distribution  $f(x/y)$  of  $X$  given  $Y$ . Patil and Rao (1974) show that the pdf of the resulting weighted version of the r.v.  $X$ , (say  $X^w$ ), is given by  $f^w(x) = \mathbb{E}[w(Y)/x]f(x)/\mathbb{E}[w(Y)]$ . In this case the weight function  $w(x, y) = w(y)$  is the regression of  $w(Y)$  on  $X$ . As is the case in sample surveys involving sampling with probability proportional to size,  $Y$  may define a design variable, and the sample is selected with probability proportional to the  $Y$  values. Thus, the original pdf is adjusted according to a given design of

selection of samples so as to improve the efficiency of estimators.

Cox (1969) discussed an idealized model in sampling textile fibers called length-biased or weighted sampling in which the choice of selection is proportional to the length of the fiber. When sampling a cross-section of a yarn, the fiber that crosses this cross-section has a length biased version of the distribution of the lengths of all the fiber in a yarn.

### 1.3.2. Relationship between Unweighted and Weighted Distributions

Cox (1969), Patil and Rao (1978), Mahfoud and Patil (1982) compared moments of unweighted distributions and weighted distributions. Gupta and Keating (1986); Gupta and Kirmani (1980) studied the relationship between unweighted and weighted distributions in the context of reliability and life testing, in particular, relationships of survival functions, failure rates and residual life functions. The relationships of the measures depend on the nature of the weight functions. We review some of the relationships assuming  $Y$  is a length-biased random variable of the original unweighted random variable  $X$ . Thus

$$E(Y) - E(X) = V(X)/E(X) \quad (1.3.1)$$

or

$$E(Y) = E(X)[ 1 + V(X)/E^2(X) ], \quad (1.3.2)$$

where  $E(X)$  denotes the expected value of  $X$  and  $V(X)$  denotes the variance of  $X$ . Hence the length-biased  $Y$  tends to record larger values of  $X$  more often.

Gupta (1981) discusses and relates the reliability measures of the original distribution to the reliability measures of the corresponding length-biased distribution. He shows that

$$Q(y) \geq S(y), \quad \lambda_g(y) \leq \lambda_f(x) \quad \text{and} \quad r_g(y) \geq r_f(x)$$



where  $Q(y)$  and  $S(y)$  are the survival functions of the length-biased distribution and unweighted distribution respectively,  $\lambda_g(y)$  and  $\lambda_f(y)$  are the hazard functions and  $r_g(y)$  and  $r_f(y)$  are the mean residual life functions of the length-biased distribution and its original distribution, respectively.

The reliability functions are defined in terms of the random variable  $X$  with expectation  $EX = \mu$ .

The survival function is defined as:

$$S(t) = P\{X > t\} = 1 - F(t) \quad \text{and the derivative } S'(x) = -f(x)$$

and the hazard function given by

$$\lambda(t) = f(t)/S(t). \quad (1.3.3)$$

The mean residual life function is defined as:

$$r(t) = E\{X - t / X > t\} = \int_t^\infty S(u)du/S(t). \quad (1.3.4)$$

The equivalence of  $S(t)$ ,  $\lambda(t)$  and  $r(t)$  are can be established as follows:

$$\begin{aligned} S(t) &= \exp\left\{-\int_0^t \lambda(u)du\right\} \\ &= \frac{r(0)}{r(t)} \exp\left\{-\int_0^t \frac{du}{r(u)}\right\}, \end{aligned} \quad (1.3.5)$$

showing that  $S(t)$  can be derived from  $\lambda(t)$  or  $r(t)$ .

The hazard function can be expressed in terms of  $r(t)$  as

$$\lambda(t) = \{1 + r'(t)\} / r(t), \quad (1.3.6)$$

and  $r$  can be obtained from  $S$  by the formula,

$$r(t) = \int_t^\infty S(u)du/S(t) = \{\mu - \int_0^t S(u)du\}/S(t). \quad (1.3.7)$$

Hence given any one of the three measures, the other two can be determined.

To aid in establishing the relation between the reliability measures of unweighted distribution and its length-biased form of distribution due to Gupta (1981), we use the notations of the reliability measures as defined above.

The relation between survival functions is given by observing that

$$Q'(y) = yS'(y)/\mu, \quad (1.3.8)$$

and integrating both sides to obtain

$$Q(y) = -\int_y^{\infty} \frac{uS'(u)}{\mu} du = S(y)\{y + r_f(y)\}/\mu. \quad (1.3.9)$$

Using the definition of  $\lambda(u)$ , the relation between the hazard functions is established below as:

$$\begin{aligned} \lambda_g(y) &= \frac{g(y)}{Q(y)} = y\lambda_f(y)/\{y + r_f(y)\} \\ &= y\lambda_f(y)/\{y + \exp\{\Lambda_f(y)\}[\mu - \int_0^y \exp\{-\Lambda(t)\}dt]\}, \end{aligned} \quad (1.3.10)$$

where  $\Lambda(t) = \int_0^t \lambda(u)du$  is the integrated hazard function.

From algebraic manipulation and using the definition of  $r(t)$ , we obtain  $\lambda_f(y)$  as a function of  $\lambda_g(y)$  given by

$$\begin{aligned} \lambda_f(y) &= \frac{f(y)}{S(y)} = \frac{g(y)/y}{\int_y^{\infty} \frac{1}{u} g(u)du} = \frac{\lambda_g(y)/y}{\int_y^{\infty} [\lambda_g(u)Q(u) \frac{1}{uQ(y)}] du} \\ &= \frac{\lambda_g(y)/y}{\int_y^{\infty} [\lambda_g(u) \frac{1}{u} \exp\{-\int_y^u \lambda_g(s)ds\}] du}. \end{aligned} \quad (1.3.11)$$

The relation between the mean residual functions is established below:

From the definition of  $r(t)$  and (1.3.9), the mean residual life function,  $r_g(y)$ , of the length-biased distribution is expressed in terms of the mean residual life function,  $r_f(y)$ , of the unweighted distribution as;

$$\begin{aligned} r_g(y) &= \int_y^{\infty} \frac{Q(u)}{Q(y)} du = \int_y^{\infty} \frac{S(u)\{u+r_f(u)\}}{S(y)\{y+r_f(y)\}} du \\ &= \frac{r_f(y)}{y+r_f(y)} \int_y^{\infty} \frac{u+r_f(u)}{r_f(u)} \exp\left\{-\int_y^u \frac{ds}{r_f(s)}\right\} du \end{aligned} \quad (1.3.12)$$

and  $r_f(y)$  is expressed in terms of  $r_g(y)$  as:

$$y + r_f(y) = \frac{\exp\left[-\int_0^y \{du/r_g(u)\}\right]/r_g(y)}{\int_0^{\infty} \{1+r_g'(t)\} / \{t(r(t))^2\} \exp\left[-\int_0^t \{du/r_g(u)\}\right] dt} . \quad (1.3.13)$$

These results show that if the data are unknowingly obtained from a weighted version of the original distribution, the analyst will overestimate the residual life function.

### 1.3.3. Comparison of Unweighted and Weighted Experiments

The theory of comparison of statistical experiments was developed originally by Blackwell (1951), based on the concept of sufficiency. For two experiments  $E_x = \{X, S_x ; f_{\theta}, \theta \in \Omega\}$  and  $E_y = \{Y, S_y ; g_{\theta}, \theta \in \Omega\}$  where  $E_x$  denotes an experiment in which a random variable defined on some sample space  $S_x$  is to be observed, and the distribution  $f_{\theta}$  of  $X$  depends on the unknown parameter  $\theta \in \Omega$ ,

and likewise definition for  $E_y$ , Blackwell defines the experiment  $E_x$  to be sufficient for experiment  $E_y$  if there exists a stochastic transformation  $Z(X)$  such that for each  $\theta \in \Omega$ , the random variable  $Z(X)$  and  $Y$  have identical distributions. This means that if  $\mathbf{x}$  is observed from performing experiment  $E_x$ , then  $\mathbf{x}'$  obtained according to distribution  $h(z/\mathbf{x})$  of  $Z(X)$  will be as informative as  $\mathbf{x}^*$  resulting from experiment  $E_y$ .

Lindley (1955) gives the comparison of two experiments,  $E_x$  and  $E_y$  with the same parameter  $\theta \in \Omega$ , based on the measure of information provided by the experiments. He defines experiment  $E_x$  to be more informative than experiment  $E_y$  if  $I_x(\theta) \geq I_y(\theta)$ , where  $I_x(\theta)$  is the Fisher information matrix from experiment  $E_x$ . If the equality holds, the two experiments are said to be equally informative. Lindley shows that if  $E_x$  is sufficient for  $E_y$  then  $E_x$  is not less informative than  $E_y$ .

Bayarri and De Groot (1989) compared information about  $\theta$  obtained from weighted distributions with that obtained from unweighted distribution, based on the concepts of sufficiency as developed by Blackwell. They defined experiment  $E_x$  to be sufficient for another experiment  $E_y$ , with the same parameter space, if  $I_x(\theta) - I_y(\theta)$  is nonnegative definite. Their results applied to the Binomial, Poisson and Negative binomial, and their respective versions of length-biased distributions show different effects for each family of distributions. For the binomial,  $E_x$  is sufficient for  $E_y$ ; for the Poisson  $E_x$  and  $E_y$  are equivalent; and for the negative binomial  $E_y$  is sufficient for  $E_x$ .

Another form of weighted sample design, selection sample, arise when sample selection is restricted to a certain subset of the sample space. The weighted form of the original pdf  $f$  is expressed as

$$g(x/\theta) = f(x/\theta)/s(\theta), \quad (1.3.14)$$

where  $s(\theta) = P(X \in S/\theta)$ . Clearly this is another form of (1.1.1) with  $w(x) = 1$ , if  $X \in S$ ;  $w(x) = 0$ , if  $X \notin S$ .

Bayarri and De Groot (1986) compared the experiment in which a selection sample is obtained with that of the unweighted distribution by means of Fisher information. They considered the case where  $X$  has a normal distribution with unknown mean  $\theta$  and known precision, and studied various selection models for the restricted sample. If  $Y$  is restricted to the set  $Y \geq \tau$ , it is shown that  $E_x$  is sufficient for  $E_y$ ; if  $Y$  is restricted to the set  $S = \{y: y \leq \tau_1 \text{ or } y \geq \tau_2\}$ , where  $\tau_1 \leq \tau_2$ , then  $E_y$  is sufficient for  $E_x$ ; and for  $\tau_1 \leq y \leq \tau_2$ , they show  $E_x$  to be sufficient for  $E_y$ . Further it is shown by Bayarri and De Groot (1986) that for selection models from an exponential family with density

$$f(x/\theta) = a(x)c(y)\exp\{v(\theta)u(x)\},$$

the information matrix for the selection model is given by

$$I_y(\theta) = I_x(\theta) + \frac{d^2}{d\theta^2} \log s(\theta),$$

where  $I_x(\theta)$  denotes the information matrix from the unrestricted sample. Hence

$$I_y(\theta) \geq I_x(\theta) \quad \text{iff } \log s(\theta) \text{ is convex, and}$$

$$I_y(\theta) \leq I_x(\theta) \quad \text{iff } \log s(\theta) \text{ is concave.}$$

#### 1.3.4. Estimation of parameters

Cox (1969) considered estimating the mean of the unweighted distribution  $f(y)$  based on the weighted pdf  $g(y)$ , and also considered estimation of the unweighted distribution function  $F$  at a fixed  $y > 0$ . Denoting by  $E_g$  the mean with respect to the weighted distribution, the relations between the moments of  $g(y)$  and the moments of  $f(y)$  can be expressed as:

$$\mathbb{E}_g(Y^r) = \frac{\mu_r + 1}{\mu} \quad (1.3.15)$$

where  $\mu_r$  is the  $r^{\text{th}}$  moment of the r.v.  $Y$  with respect to pdf  $f$ .

From (1.3.15) we obtain

$$\mathbb{E}_g(Y) = \mu \left( 1 + \frac{\sigma^2}{\mu^2} \right), \quad (1.3.16)$$

where  $\mu$  and  $\sigma^2$  are the mean and variance of the unweighted distributions.

Also from (1.3.15) we get

$$\mathbb{E}_g(Y^{-1}) = \frac{1}{\bar{\mu}} \quad \text{and} \quad \text{Var}(Y^{-1}) = \frac{1}{\mu^2} [\mu\mu_{-1} - 1],$$

which suggests using

$$\frac{1}{\bar{\mu}} = \frac{1}{n} \left\{ \sum Y_i^{-1} \right\}, \quad (1.3.17)$$

as an unbiased estimate of  $\frac{1}{\mu}$ . Cox (1969) proposed  $\tilde{\mu}$ , the harmonic mean of the values of  $y$  in the biased sample, as the estimator of the mean of the original distribution.

As discussed by Sen (1987),  $\tilde{\mu}$  is a biased though consistent and asymptotically normal estimator of  $\mu = \mathbb{E}_f(Y)$  with the bias of order  $n^{-1}$ , since

$$\mathbb{E}[\tilde{\mu}] = \mu \{ 1 + n^{-1}(\mu\mu_{-1} - 1) \} + o\left(\frac{1}{n}\right), \quad (1.3.18)$$

showing the leading bias in the estimator  $\tilde{\mu}$  as  $n^{-1}\mu(\mu\mu_{-1} - 1)$ . The bias reduction in the estimator of  $\mu$  can be achieved by jackknifing, as suggested by Sen (1987). It is of interest to investigate the properties of the estimate of the mean  $\mu$  for finite samples when the mean is considered a linear function of some covariate terms.

Krieger and Pfefferman (1992), considered approaches for maximum likelihood estimation from survey data. The two approaches, the classical pseudo likelihood (CPL) and MLE derived from weighted distributions (WDML) were compared for their performance using a simulation study.

The CPL utilizes the sample selection probabilities to estimate the census likelihood equations. Suppose the population values  $Y_i$  are independent with a common distribution  $f(Y; \theta)$  and  $L_N(\theta; Y) = \sum_{i=1}^N \log f(Y_i; \theta)$  define the census log-likelihood function. The MLE,  $\hat{\theta}$  is a solution to the equation,

$$\tilde{U}(\theta) = \partial L_N(\theta; Y) / \partial \theta = \sum_{i=1}^N u(\theta; y_i) = 0 \quad (1.3.19)$$

where  $u(\theta; y_i) = \frac{\partial}{\partial \theta} \log f(Y_i; \theta)$ .

Assume only the sample values  $\{y_i, i \in s\}$  and the sample selection probabilities  $\{P_i, i \in s\}$  are known. The pseudo MLE uses the estimator  $\hat{U}(\theta)$  of  $\tilde{U}(\theta)$  and is thus the solution to

$$\hat{U}(\theta) = \sum_{i=1}^n w_i^* u(\theta; y_i), \quad (1.3.20)$$

where for selection without replacement  $w_i^* = [1/P_i (i \in s)]$ , and for selection with replacement  $w_i^* = (1/nP)$ .

The second approach WDML is based on the concept of weighted distributions and utilizes the sample selection probabilities to adjust the original unweighted distribution. The MLE of  $\theta$  is then the solution to the maximum likelihood estimating equations of the weighted distribution. Krieger and Pfefferman illustrate the comparison with an example in which  $X'_i = (Y'_i, Z'_i)$  are

independent realizations from the multivariate normal distribution with mean  $\underline{\mu}_X = (\underline{\mu}'_Y, \underline{\mu}'_Z)$  and variance-covariance, V-C, matrix

$$\Sigma_{XX} = \begin{bmatrix} \Sigma_{YY} & \Sigma_{YZ} \\ \Sigma_{ZY} & \Sigma_{ZZ} \end{bmatrix}.$$

The results obtained when estimating  $\mu_Y = E(Y)$  and  $\sigma^2 = \text{Var}(Y)$  were compared for the two distributions for their Residual Mean Square Error (RMSE). It is shown from a simulation study that WDML dominates CPL. Krieger and Pfefferman state as a possible explanation for this result the fact that WDML is more “model dependent” whereas the CPL is viewed as “estimator maximizing the design unbiased estimator of the likelihood equations holding for the population”.

Godambe and Rajarshi (1989) extended the idea of optimal estimating equations for semi-parametric models to weighted distributions for a finite population of size N. Assuming observational samples  $y_i$  are made with probability  $w(y_i)$  they considered the maximum likelihood equations for  $\theta$  based on the sample data  $d$ , from a finite population  $\mathcal{P}$

$$d = \{ (i, y_i) : i \in s \}, \quad s \subset \mathcal{P}.$$

The probability density of  $d$  is given by

$$\text{Prob}(d; \theta) = \prod_{i \in s} \frac{w(y_i) f(y_i; \theta)}{\mu(\theta)} \prod_{i \in s} \mu(\theta) \prod_{i \notin s} (1 - \mu(\theta)), \quad (1.3.21)$$

where  $\mu(\theta) = \int w(y_i) f(y_i; \theta) dy_i$ , for  $i = 1, 2, \dots, N$ .

The maximum likelihood estimator for  $\theta$  is then found in the usual manner as



a solution to the maximum likelihood equations given by

$$\sum_{i \in s} \frac{\partial}{\partial \theta_j} \log \frac{w(y_i) f(y_i; \theta)}{\mu(\theta)} + \frac{n - \mu(\theta)N}{\mu(\theta)(1 - \mu(\theta))} \frac{\partial \mu(\theta)}{\partial \theta_j} = 0, \quad j = 1, 2, \dots, k, \quad (1.3.22)$$

where  $n$  is the sample size.

The maximum likelihood equations in (1.3.22) include an additional term

$$\frac{n - \mu(\theta)N}{\mu(\theta)(1 - \mu(\theta))} \frac{\partial \mu(\theta)}{\partial \theta_j} \quad (1.3.23)$$

to the usual ML equations, given by the first term on the r.h.s. of (1.3.22). The population size  $N$  can be treated as a parameter to be jointly estimated with  $\theta$  from the distribution (1.3.21).

The theory of estimating functions as discussed by Godambe and Rajarshi provides equations that are jointly optimal for  $\theta$  and  $N$ . They considered an experiment in which each individual  $i \in \mathcal{P}$  has a chance of selection  $m(\theta)$ , and  $y_i$  is drawn from the distribution  $g(y; \theta)$  in (1.2.1), such that the sample data consists of  $(i, y_i)$ . To estimate  $\theta$ , the  $2N$  estimating functions are constructed for each  $i \in \mathcal{P}$  and  $j$ , as

$$h_i = \phi_i \frac{\partial}{\partial \theta_j} \log \frac{w(y_i) f(y_i; \theta)}{\mu(\theta)} \quad \text{and} \quad h_{i+N} = \phi_i - \mu(\theta), \quad i = 1, 2, \dots, N, \quad (1.3.24)$$

where  $\phi = 1$  if the  $i^{\text{th}}$  individual is selected in the sample and 0 otherwise. Since  $\mathbb{E}_g(h_i) = 0$  and  $\mathbb{E}_g(h_i h_{i'}) = 0$ , for  $i, i' = 1, 2, \dots, N$ , the  $2N$  estimating functions are mutually orthogonal. Then according to Godambe and Rajarshi, the optimal estimating combination of the  $h_i$ 's is given by

$$\begin{aligned}
h &= \sum_{i=1}^N \left\{ h_i \frac{E_g \partial h_i / \partial \theta_j}{E(h_i^2)} + h_{i+N} \frac{E_g \partial h_{i+N} / \partial \theta_j}{E(h_{i+N}^2)} \right\} \\
&= \sum_{i=1}^N \phi_i \frac{\partial}{\partial \theta_j} \log \frac{w(y_i) f(y_i; \theta)}{\mu(\theta)} + \sum_{i=1}^N (\phi_i - \mu(\theta)) \frac{\partial \mu(\theta) / \partial \theta_j}{\mu(\theta)(1 - \mu(\theta))}, \quad (1.3.25)
\end{aligned}$$

which reduces to equation (1.3.22). Godambe and Rajarshi suggest that if  $N$  is unknown, the functions  $h_{i+N}$ ,  $i = 1, 2, \dots, N$  be combined into a single function  $\sum_{i=1}^N h_{i+N} = n - \mu(\theta)N$ , which is orthogonal to the  $h_i$ 's. Hence equations  $h = 0$  and  $n - \mu(\theta)N = 0$  provide joint optimal estimation for  $\theta$  and  $N$ .

Replacing the distributional assumption for the  $y_i$ 's by assuming only the mean  $\mu$ , variance  $\sigma^2$ , and coefficient of skewness  $\gamma$ , Godambe and Rajarshi extends the idea of estimating equations to the semi-parametric models as follows:

Let

$$E_g[Y_i] = \frac{1}{\mu}(\sigma^2 + \mu^2) = c(\theta), \quad (1.3.26)$$

$$E_g[Y_i - c(\theta)]^2 = \sigma^2 \left\{ \left( \frac{\gamma \sigma}{\mu} \right) + 1 - \left( \frac{\sigma^2}{\mu^2} \right) \right\}. \quad (1.3.27)$$

For  $j = 1, \dots, k$ , define the mutually orthogonal estimating equation as

$$h_i = \phi_i(y_i - c(\theta)) \quad \text{and} \quad h_{i+N} = \phi_i - \mu(\theta), \quad \text{for } i=1, 2, \dots, N \quad (1.3.28)$$

Then the optimum combination of the estimating functions is given by

$$h = \sum_{i=1}^N \phi_i(y_i - c(\theta)) \frac{\partial c(\theta) / \partial \theta_j}{E_g[y_i - c(\theta)]^2} + \sum_{i=1}^N (\phi_i - \mu(\theta)) \frac{\partial \mu(\theta) / \partial \theta_j}{\mu(\theta)(1 - \mu(\theta))},$$

$$\text{for } j = 1, \dots, k \quad (1.3.29)$$

where  $\phi_i = 1$  if  $y_i$  is in the sample and  $\phi_i = 0$  otherwise,  $c(\theta)$  and  $E_g[y_i - c(\theta)]^2$  are

given by (1.3.26) and (1.3.27), respectively,

and

$$\frac{\partial \mu(\boldsymbol{\theta})}{\partial \theta_j} = - \left\{ \frac{1}{\mu(\boldsymbol{\theta})} \frac{\partial \sigma^2(\boldsymbol{\theta})}{\partial \theta_j} - \left[ \frac{\sigma^2(\boldsymbol{\theta})}{\mu^2(\boldsymbol{\theta})} - 1 \right] \frac{\partial \mu(\boldsymbol{\theta})}{\partial \theta_j} \right\} \quad (1.3.30)$$

which gives

$$\begin{aligned} h = \sum_{i \in s} (y_i - c(\boldsymbol{\theta})) & \frac{\{(\sigma/\mu)^2 - 1\}(\partial \mu / \partial \theta_j) - (\partial \sigma^2 / \partial \theta_j) / \mu}{\sigma^2 \{(\gamma \sigma / \mu) + 1 - (\sigma / \mu)^2\}} \\ & + \sum_{i \in s} (\partial \log \mu(\boldsymbol{\theta}) / \partial \theta_j) + \sum_{i \notin s} \{\partial \log (1 - \mu(\boldsymbol{\theta})) / \partial \theta_j\}. \end{aligned} \quad (1.3.31)$$

The semi-parametric estimator of  $\boldsymbol{\theta}$  is obtained from (1.3.31) by solving the optimal estimating equation  $h = 0$  for  $\boldsymbol{\theta}$ .

### 1.3.5 Density Estimation

Most research has been devoted to estimating the probability distribution function  $F$  of  $f$  and the mean of the distribution based on the length-biased samples. Winter and Foldes (1988) proposed a nonparametric estimator of the cumulative hazard function and used it to construct an estimator of  $F$ . They considered a product-limit estimator for the life-testing situation where the life of  $n$  objects observed has the distribution  $G$ , the length-biased form of  $F$ , and established strong uniform consistency results for their estimator. Bhattacharyya, Franklin and Richardson (1988) compared the mean square errors of the density estimators of the weighted and unweighted distributions. Assuming  $X_1, \dots, X_n$  are i.i.d. random variables with density  $g(x) = xf(x)/\mu$ , for  $x \geq 0$  they considered an estimator  $\hat{g}(x)$ , the kernel estimator as discussed by Parzen (1962). The estimator

is of the form

$$f_n(x) = \frac{1}{nh_n} \sum K\{(x - X_j)/h_n\}, \quad (1.3.32)$$

where  $\{h_n\}$  is a sequence of constants converging to zero, and  $K(y)$  is a bounded integrable function satisfying the conditions  $\lim |yK(y)| = 0$ .

Parzen gives the asymptotic variance of this estimator, under precise conditions assumed concerning the function  $K$ , and show asymptotic normality. The motivation of this estimator is that a common estimator of the distribution function  $F(x)$  is the empirical distribution function defined as

$$F_n(x) = \frac{1}{n} \sum u(x - X_j), \text{ where } u(y) = 1 \text{ if } y \geq 0 \text{ and } u(y) = 0 \text{ otherwise.}$$

Bhattacharyya (1988) proposed estimating  $\hat{f}(x) = \hat{\mu} \times \hat{g}(x)/x$ , where  $\hat{\mu}$  is an estimator of  $\mu$ ,  $\hat{\mu} = n[\sum x_i^{-1}]^{-1}$ . It was shown that the mean square error of the length biased estimator approaches zero as  $n$  increases to infinity.

Jones (1991) derived a direct kernel density estimator for the length-biased data of the form

$$\hat{f}_n(x) = n^{-1} \hat{\mu} \sum X_j^{-1} K_h(x - X_j), \quad (1.3.33)$$

where  $h$  is the smoothing parameter and  $K_h = h^{-1}K(x/h)$ . The measure of error used for this estimator is the mean integrated squared error (IMSE) defined as

$$E \int [\hat{f}_n(x) - f(x)]^2 dx. \quad (1.3.34)$$

Jones compared the integrated mean squared errors of the estimators based

on the weighted and unweighted models. Using simulated sample data from the density  $f$  and generated length-biased data from  $g$  the results, (Jones 1991) show the estimator based on length-biased data is inferior to the estimator based on the unweighted sample in terms of integrated mean square error, IMSE. Also the  $\text{IMSE}(\hat{f})$  is less than the IMSE of the estimator provided by Bhattacharyya.

### 1.3.6 Discussion and Summary of Literature Review

Work on length-biased data analysis has focused mainly on estimation of the mean  $\mu$  and the distribution function  $F$  of the original distribution based on length-biased data. The length-biased density function considered in the literature is of the form  $g(x) = x \times f(x) / \mu$ ,  $x \geq 0$  and  $g(x) = 0$ ,  $x \leq 0$ , where  $\mu$  is the mean under the unweighted distribution.

A natural consistent estimator of  $\mu$  is given by  $\tilde{\mu} = n[\sum x^{-1}]^{-1}$ , Cox (1969) and Sen (1987). The estimators for  $\mu$  and  $F$  from the length-biased data from  $g(x)$  are compared with those obtained from the original unweighted data from  $f(x)$  by finding their asymptotic variances. The estimators are derived in a nonparametric setup, and hence do not take advantage of the distributional information of the data. Furthermore, the mean is not considered dependent on a set of possible explanatory variables which might be available in research situations.

Godambe and Rajarshi considered a parametric model and semi-parametric models for the length-biased density for a finite population of size  $N$ . Their work is limited to estimating the parameters of the length-biased distribution and does not consider estimation of the functions of the parameters and their properties. Moreover, it is not clear how the optimality of the estimating functions they

consider translates to the optimality of the parameters. It is also not clear how the number of the arbitrary estimating functions,  $h_i$ 's, is to be determined, especially in an infinite population. No attempt is made to compare their results with estimates that would be obtained from under the assumed original unweighted model.

It is the purpose of this research to extend the estimation to the situation where  $\mu$  depends on a set of covariates in a regression setup, and to investigate the properties of such estimators. The two sampling designs, sampling under  $f(x)$  and sampling under  $g(x)$  will be compared for their lack of information about the parameters by way of optimality criteria using Fisher information matrix. The next section discusses the proposed research on length-biased distributions.

## 1.4 Research Proposal

The purpose of this research is to obtain estimates of the parameters of the original (unweighted) distribution from the data generated by a length-biased design. We study the large sample properties of the maximum likelihood and the jackknifed estimators for the two designs. Both the MLE and the jackknifed estimators (originally proposed by Quenoulli) are shown to have equivalent asymptotic normal distribution. The reduction in the bias of the MLE achieved by the jackknife estimator is provided. The MLE and the jackknife estimators of the regression model are also studied for their large sample properties, both the bias and the asymptotic normal distributions.

Our work differs from that in the existing literature in that we assume a parametric estimation of the parameters. Whereas attention in the existing literature emphasizes estimating a single parameter  $\mu$ , the mean of the random

variate, from a length-biased distribution, we consider estimating a  $k$ -parameter vector  $\theta$  and for the regression model. We assume  $\mu$  depends on a set of nonstochastic covariates whose vector of coefficients  $\beta$  is the primary objective of the estimation. Furthermore, this research aims at comparing the information about the parameter estimates from the two designs by way of the Fisher information.

We concern ourselves with a random variate  $Y$ , following an exponential probability distribution with pdf given by;

$$f(y;\theta) = \exp\left\{ \sum_{j=1}^k \theta_j T_j(y) + \Psi(\theta) + C(y) \right\}, \quad (1.4.1)$$

where  $\theta$  is a  $k$ -dimensional vector of unknown parameters. Since  $f(y;\theta)$  is of exponential family type of distributions, its length-biased version  $g(y;\theta)$  also belongs to the exponential family of distributions and is given by,

$$g(y;\theta) = \exp\left\{ \sum_{j=1}^k \theta_j T_j(y) + \Psi^*(\theta) + C^*(y) \right\}, \quad (1.4.2)$$

where  $\Psi^*(\theta) = \Psi(\theta) - \log \mu(\theta)$ ,  $\mu(\theta) = E[Y]$  and  $C^*(y) = C(y) + \log(y)$ .

Our purpose is to consider estimation of the parameters or functions of the parameters  $\theta$  of the model (1.4.1) from the data generated by the probability distribution model (1.4.2). The asymptotic bias and distribution of such estimators is also established. The performance of the estimators is measured by the root mean squared error (RMSE), which is a combination of bias and variance of the estimator.

This dissertation is organized as follows:

Chapter 2 discusses the large sample properties of the MLE and the jackknife

estimator of the MLE from the data following the probability distribution (1.4.1). The jackknife variance estimator is shown to be consistent estimator of the variance of the MLE. In particular, the proposed jackknife procedure is applied to the maximum likelihood estimator  $\hat{\theta}$  of  $\theta$ . The jackknife method of estimation generates replicates of the statistic  $\hat{\theta}$  by deleting one observation from the sample and then computes the statistic from the remaining data. The bias of the jackknife estimator is compared with the bias of the MLE. It is shown that the jackknife estimators have the same asymptotic normal distribution as the MLE's.

Chapter 3 extends the results of Chapter 2 to the case of estimating the parameters of the regression model. Regression analysis concerns the study of the relationship between the response variable  $Y$  and a set of predictor variables  $\mathbf{Z} = (Z_1, \dots, Z_p)^T$ . Given that the data follows an exponential family of distributions models, we hope to obtain the small sample results for the estimators, failing which we shall concern ourselves only with asymptotic results.

The model considered is given by

$$E[Y/Z] = \mu = \mu(\boldsymbol{\beta}, \phi), \quad (1.4.3)$$

where  $\boldsymbol{\beta}$  is a  $p$ -vector parameter of regression coefficients and  $\phi$  is the scale parameter. The  $\mathbf{Z}$  matrix is nonstochastic and is assumed to be known. The parameter of interest  $\hat{\boldsymbol{\beta}}$  of  $\boldsymbol{\beta}$  is estimated by maximum likelihood estimation methods. By deleting each row of  $\mathbf{Y}$  and  $\mathbf{Z}$ , the jackknife estimator is derived. We study the large sample distribution of the MLE and the jackknife estimators.

Chapter 4 discusses the MLE and the jackknife estimators under the length-biased distributions. The bias and large sample properties of the estimators are derived. We compare the Fisher information matrix of the parameters of the



original distribution with the Fisher information matrix from the length-biased distribution, by way of D-optimality and E-optimality (Keifer, 1975).

Chapter 5 compares the jackknife estimates and the MLE computed from (1) the 1990 Sudan DHS birth interval data, and (2) the simulated data. This is intended to demonstrate the applications of our methodology. Two distributions considered in the simulation study for the exponential family of distributions are the lognormal and the gamma distributions. Simulation studies are undertaken on the data generated from these distributions to illustrate the computations of jackknife estimators. These studies will indicate how well our estimators perform in small to moderate sample sizes. Attention is centered on estimating and making statistical inferences about the parameters of the original distributions based on the length-biased data. Hence the bias and variance estimation is of significant importance in this study. Chapter 6 summarizes the results of the dissertation and provides suggestions for future research.

# Chapter 2.

## Parametric Estimation From The Exponential Family.

### 2.1 Introduction.

Our objective in this chapter is to derive expressions for the asymptotic bias of the MLE and the Jackknife Estimators of the parameters  $\theta$  of the Exponential family of distributions. We also study the asymptotic distribution of the estimators and show the MLE and Jackknife estimators of the parameters of the Exponential family of distributions are asymptotically equivalent, in the sense that  $\sqrt{n}(\hat{\theta}_n - \theta)$  and  $\sqrt{n}(\tilde{\theta}^J - \theta)$  have identical asymptotic normal distribution, where  $\hat{\theta}$  and  $\tilde{\theta}^J$  are the MLE and jackknifed estimators respectively. The convergence of the sample variance of the pseudo-values (Tukey) to the variance of the MLE is shown.

### 2.2 Maximum Likelihoods Estimators

We consider the Maximum Likelihood Estimators (MLE) of parameters of the Exponential family of distributions. In this section we show the conditions

under which the MLE of the parameters of the Exponential family are consistent. The asymptotic bias and distribution of the MLE from the exponential family of distributions are derived.

Suppose that  $X_1, \dots, X_n$  are independent and identically distributed (i.i.d.) random variables with density function  $f(X; \theta)$ ,  $\theta \in \Theta \subset \mathbb{R}^k$ . For the Exponential family considered in this section, the density is  $f(X; \theta)$  the form

$$f(X; \theta) = c(X) \exp \left\{ \sum_{j=1}^k \theta_j T_j(X) + \Psi(\theta) \right\} \quad (2.2.1)$$

where  $c(\cdot)$  and  $\Psi(\cdot)$  are some known functions. The likelihood function for the sample  $X_1, \dots, X_n$  of size  $n$  is defined by

$$L_n(\theta; \mathbf{X}) = \prod_{i=1}^n f(X_i; \theta). \quad (2.2.2)$$

Whatever its value, it is well known that the likelihood function is maximized at the true value  $\theta \in \Theta$ . Thus the principle of maximum likelihood estimation consists of choosing an estimate of  $\theta$ , say  $\hat{\theta}$  that maximizes (2.2.2). It is often convenient to work with the log-likelihood function and since the log function is monotone, the value of  $\hat{\theta}$  that maximizes  $\log L_n(\theta; \mathbf{X})$  will also maximize  $L_n(\theta; \mathbf{X})$ . Thus the MLE of  $\theta$  may be obtained by maximizing the sample log-likelihood function and in which case, is the solution to the likelihood equation,

$$U_n(\theta) = \frac{\partial}{\partial \theta} \log L_n(\theta; \mathbf{X}) = \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(X_i; \theta) = \mathbf{0}. \quad (2.2.3)$$

The usual assumptions made for the MLE to exist and have asymptotic

normality properties are:

A1.  $\frac{\partial}{\partial \theta_j} \log f(X_i; \theta)$  and  $\frac{\partial}{\partial \theta_j \partial \theta_l} \log f(X_i; \theta)$  for  $j, l = 1, \dots, k$  exist almost everywhere, and

$$\left| \frac{\partial}{\partial \theta_j} \log f(X_i; \theta) \right| \leq K_1(X), \quad \left| \frac{\partial^2}{\partial \theta_j \partial \theta_l} \log f(X_i; \theta) \right| \leq K_2(X)$$

where  $K_1(X)$  and  $K_2(X)$  are integrable functions, and  $\int K_i(X) < \infty$  ( $i = 1, 2$ ).

A2.  $0 < I(\theta) = E_{\theta} \left\{ \frac{\partial}{\partial \theta} \log f(X; \theta) \right\} \left\{ \frac{\partial}{\partial \theta} \log f(X; \theta) \right\}^T < \infty$

$I(\theta)$  is called the Fisher information matrix.

A3.  $E_{\theta} \left\{ \sup_{|h| < \delta} \left\| \frac{\partial^2}{\partial \theta \partial \theta^T} \log f(X_i; \theta) - \frac{\partial^2}{\partial \theta \partial \theta^T} \log f(X_i; \theta + h) \right\| \right\} \longrightarrow 0$ , as  $\delta \rightarrow 0$ .

The implications of the above set of assumptions for the particular exponential family of distributions (2.2.1) studied here lead to the following assumptions :

B1.  $\Psi(\theta) = \Psi(\theta_1, \dots, \theta_n)$  has continuous partial derivatives at all  $\theta_j$   $j = 1, \dots, k$

B2.  $\left( - \frac{\partial^2}{\partial \theta_j \partial \theta_l} \Psi(\theta) \right)$ ,  $j, l = 1, \dots, k$ , is a  $(k \times k)$  positive definite matrix for all  $\theta \in \Theta$ .

B3.  $E_{\theta} | T_j(X) | < \infty$ ,  $j = 1, \dots, k$   $\theta \in \Theta$ .

To construct the estimating equations for  $\theta$  we consider the individual contribution of each individual data point. The partial derivatives of the log likelihood of (2.2.1) with respect to  $\theta$  as,

for  $j, l = 1, \dots, k$

$$\frac{\partial}{\partial \theta_j} \log f(X; \theta) = T_j(X_i) + \frac{\partial}{\partial \theta_j} \Psi(\theta) \quad (2.2.4)$$

and

$$\frac{\partial^2}{\partial \theta_j \partial \theta_l} \log f(X; \theta) = \frac{\partial^2}{\partial \theta_j \partial \theta_l} \Psi(\theta). \quad (2.2.5)$$

The MLE  $\hat{\theta}_n$ , of  $\theta$  is a solution to the MLE equations (2.2.3), so that for the probability density  $f(X; \theta)$  given by (2.2.1),

$$U_n(\hat{\theta}_n) = \sum_{i=1}^n \tilde{T}(X_i) + n \psi(\hat{\theta}_n) = 0, \quad (2.2.6)$$

which implies that

$$\frac{1}{n} \sum_{i=1}^n \tilde{T}(X_i) + \psi(\hat{\theta}_n) = 0, \quad (2.2.7)$$

where  $\psi(\hat{\theta}_n) = \frac{\partial}{\partial \theta} \Psi(\theta) \Big|_{\hat{\theta}_n}$  and  $\tilde{T}(X_i) = (T_1(X_i), \dots, T_k(X_i))^T$ .

Since

$$0 = E_{\theta} \frac{\partial}{\partial \theta} \log f(X_i; \theta) = \tilde{T}(X_i) + \psi(\theta), \quad (2.2.8)$$

we have

$$E_{\boldsymbol{\theta}} \mathbb{T}(X_i) = -\boldsymbol{\psi}(\boldsymbol{\theta}). \quad (2.2.9)$$

The first term on the LHS of (2.2.7) is the mean of a sample of size  $n$  from a population having mean  $-\boldsymbol{\psi}(\boldsymbol{\theta})$ , so by Khintchine's SLLN, as  $n \rightarrow \infty$ ,

$$\frac{1}{n} \sum_{i=1}^n \mathbb{T}(X_i) \xrightarrow{\text{a.s.}} -\boldsymbol{\psi}(\boldsymbol{\theta}). \quad (2.2.10)$$

By (2.2.3)

$$\frac{\partial U'_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^T} = \sum_{i=1}^n \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \log f(X_i; \boldsymbol{\theta}) = n \boldsymbol{\varphi}(\boldsymbol{\theta}) = n \left( \frac{\partial^2}{\partial \theta_j \partial \theta_l} \Psi(\boldsymbol{\theta}) \right) \quad (2.2.11)$$

and differentiating (2.2.1) twice with respect to  $\boldsymbol{\theta}$ , under the integral sign, and using B2, we have

$$-\boldsymbol{\varphi}(\boldsymbol{\theta}) = \left( -\frac{\partial^2}{\partial \theta_j \partial \theta_l} \Psi(\boldsymbol{\theta}) \right) = \text{Var}_{\boldsymbol{\theta}}(\mathbb{T}(X)) > 0, \quad (2.2.12)$$

we see that the log-likelihood is concave. Therefore if a solution to the likelihood equation (2.2.3) exists, it must be an MLE.

By assumption B1,  $\boldsymbol{\psi}(\boldsymbol{\theta}^*)$  is continuous in a neighborhood of the true parameter value  $\boldsymbol{\theta} \in \Theta$ , i.e. for every  $\epsilon > 0$  there exist  $\eta > 0$  such that

$$\| \boldsymbol{\psi}(\boldsymbol{\theta}^*) - \boldsymbol{\psi}(\boldsymbol{\theta}) \| < \epsilon, \quad \text{for all } \| \boldsymbol{\theta}^* - \boldsymbol{\theta} \| < \eta. \quad (2.2.13)$$

If  $-\psi(\theta) = E_{\theta}[T(X)]$  is a one-to-one function of  $\theta$ , then (2.2.13) necessarily hold.

**Example 1:** Let  $X_1, \dots, X_n$  be i.i.d. r.v.'s with a lognormal  $(\mu, \sigma^2)$  distribution. The likelihood function is

$$L_n(\mu, \sigma^2) = (2\pi\sigma)^{-n/2} \left( \prod_{i=1}^n X_i \right)^{-1} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (\log X_i - \mu)^2\right\},$$

Then  $\theta_1 = \frac{\mu}{\sigma^2}$ ,  $\theta_2 = -\frac{1}{2\sigma^2}$ ,  $T_1(X_1) = \log X_1$ ,  $T_2(X_1) = (\log X_1)^2$  and

$$\Psi(\theta) = \frac{1}{2} \left( \frac{\theta_1^2}{2\theta_2} - \log \frac{-1}{2\theta_2} \right), \text{ so that we have}$$

$$E_{\theta} T(X_i) = -\psi(\theta) = \frac{1}{2} \left( -\frac{\theta_1}{\theta_2}, \frac{\theta_1^2}{2\theta_2^2} - \frac{1}{\theta_2} \right)^T = (\mu, \mu^2 + \sigma^2)^T,$$

which is a one-to-one function of  $\theta$  and hence (2.2.13) is satisfied.

**Example 2 :** Let  $X_1, \dots, X_n$  be i.i.d. r.v.'s following a Poisson( $\lambda$ ) distribution. The likelihood of  $\theta$  is  $L_n(\lambda) = \exp\{-\lambda + X \log \lambda - \log X!\}$ . So that  $T(X) = X$ ,  $\theta = \log \lambda$ .

$E_{\theta} T(X) = -\psi(\theta) = \lambda = \exp(\theta)$  which is a one-to-one monotone function of  $\theta$  satisfying (2.2.13).

There are situations when the likelihood equation (2.2.3) has at least one consistent solution, so that the likelihood function has a relative maxima at such solutions. When the likelihood has multiple roots we will require that

$$\| \psi(\theta^{**}) - \psi(\theta) \| > 0, \quad \text{if } \| \theta^{**} - \theta \| > c \quad (2.2.14)$$

where  $\theta^{**}$  is outside the  $\eta$  neighborhood of the true parameter  $\theta$  and  $c > 0$ . That is, we can choose  $c$  sufficiently small in the neighborhood of  $\theta^{**}$  such that (2.2.14) is satisfied. This requirement ensures the likelihood evaluated at  $\theta$  is larger than at any other  $\theta^{**} \in \Theta$ .

Now since the first term on the LHS of (2.2.7) converges a.s. to  $-\psi'(\theta)$  and  $U_n(\hat{\theta}) = 0$ ,

$$\psi(\hat{\theta}) \xrightarrow{\text{a.s.}} \psi(\theta) \tag{2.2.15}$$

By B2 and (2.2.13) we have  $\hat{\theta} \xrightarrow{\text{a.s.}} \theta$ .

## 2.2.1 Asymptotic Distribution and Bias of the MLE

The asymptotic distribution of MLE have been studied extensively in the literature, (Kendal and Stuart, 1958; Cramer, 1946; LeCam, 1957; Wilks, 1963 and others). In this study we apply the theory of asymptotic properties of the MLE specifically to the exponential family (2.2.1), furthermore the form of bias of the MLE  $\hat{\theta}_n$  is presented to justify the jackknife estimation to reduce bias of the MLE in large samples. Our results relies heavily on Taylor expansion of the gradient of the log-likelihood function in some neighborhood of the true parameter  $\theta$ .

To establish the asymptotic normality of  $\hat{\theta}_n$  we let



$$U_n(\boldsymbol{\theta}) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial}{\partial \boldsymbol{\theta}} \log f(X_i; \boldsymbol{\theta}) \quad (2.2.16)$$

using (2.2.8), assumption B2 and the Central Limit Theorem on  $\partial/\partial \boldsymbol{\theta} \log f(X_i; \boldsymbol{\theta})$ , ( $i = 1, 2, \dots, n$ ) we note that  $U_n(\boldsymbol{\theta}) \xrightarrow{\mathcal{D}} \mathcal{N}(\mathbf{0}, I(\boldsymbol{\theta}))$ , where  $I(\boldsymbol{\theta})$  is the Fisher information matrix.

Expanding  $U_n(\hat{\boldsymbol{\theta}}_n)$  in the neighborhood of  $\boldsymbol{\theta}$  we write

$$\mathbf{0} = U_n(\hat{\boldsymbol{\theta}}_n) - U_n(\boldsymbol{\theta}) = U_n'(\boldsymbol{\theta})|_{\boldsymbol{\theta}^*}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}), \quad (2.2.17)$$

where  $\boldsymbol{\theta}^*$  is a point on the line segment connecting  $\hat{\boldsymbol{\theta}}_n$  and  $\boldsymbol{\theta}$ . Since from (2.2.13) and (2.2.15),  $\hat{\boldsymbol{\theta}}_n$  is strongly consistent for  $\boldsymbol{\theta}$  and for the exponential family of distributions (2.2.1), and with the assumption that  $\Psi(\boldsymbol{\theta})$  are continuously differentiable we have,

$$\boldsymbol{\varphi}(\boldsymbol{\theta}^*) = \boldsymbol{\varphi}(\boldsymbol{\theta}) + o_p\left(\frac{1}{\sqrt{n}}\right), \quad \text{as } n \longrightarrow \infty. \quad (2.2.18)$$

Hence we write (2.2.17) as

$$\begin{aligned} \boldsymbol{\psi}(\hat{\boldsymbol{\theta}}_n) - \boldsymbol{\psi}(\boldsymbol{\theta}) &= \boldsymbol{\varphi}(\boldsymbol{\theta})|_{\boldsymbol{\theta}^*}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) \\ &= \boldsymbol{\varphi}(\boldsymbol{\theta})(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) + o_p\left(\frac{1}{\sqrt{n}}\right), \end{aligned} \quad (2.2.19)$$

where the last equality in (2.2.19) is by virtue of (2.2.18) and by using the fact that  $(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta})$  is  $O(1)$ . Now by assumption B2,  $-\boldsymbol{\varphi}(\boldsymbol{\theta})$  is positive definite. From (2.2.19) we get

$$\sqrt{n}(\hat{\theta}_n - \theta) = \sqrt{n} [\varphi(\theta)]^{-1} [\psi(\hat{\theta}_n) - \psi(\theta)] + o_p(1), \quad \text{as } n \longrightarrow \infty. \quad (2.2.20)$$

So that by Slutsky's Theorem, the LHS of (2.2.20) converges in distribution to the same asymptotic distribution of  $\sqrt{n} [\varphi(\theta)]^{-1} [\psi(\hat{\theta}_n) - \psi(\theta)]$ . Now the first term on the RHS of (2.2.20) is

$$I^{-1}(\theta) U_n(\theta) \longrightarrow \mathfrak{D} \mathcal{N}(\mathbf{0}, I^{-1}(\theta)), \quad \text{as } n \longrightarrow \infty \text{ and hence}$$

$$\sqrt{n}(\hat{\theta}_n - \theta) \longrightarrow \mathfrak{D} \mathcal{N}(\mathbf{0}, I^{-1}(\theta)), \quad (2.2.21)$$

thus establishing the asymptotic normality of the MLE  $\hat{\theta}_n$  of  $\theta$  for the exponential family (2.2.1). The result (2.2.21) for general distributions is classical and is contained in Cramer, 1946 and Kendal and Stuart, 1958 among other books .

Although the MLE generally possess nice asymptotic properties of consistency and asymptotic normality, they are often not unbiased estimators. We wish to show the form of the bias in the MLE and later show how the jackknifing method can reduce the bias. To study the bias of  $\hat{\theta}_n$ , note that  $-\psi(\theta) = E_{\theta}[T(X_i)]$  is a one-to-one function of  $\theta$  (eg. if  $X$  follows a Poisson distribution with mean  $\lambda$ , then  $E[T(X)] = -\psi(\theta) = e^{\theta} = \lambda$ , so that there exists a transformation  $\mathbf{H} = \mathbf{H}(-\psi(\theta))$  satisfying

$$\mathbf{H}(-\psi(\theta)) = \theta. \quad (2.2.22)$$

**Proposition 1.** Assume the set of assumptions "B". Define  $\eta = -\psi(\theta)$  and a

function  $\mathbf{H}: \mathbb{R}^k \rightarrow \mathbb{R}^k$  given by  $\mathbf{H}(\boldsymbol{\eta}) = \boldsymbol{\theta}$ . Assume  $\mathbf{H}(\boldsymbol{\eta})$  is continuously differentiable in some neighborhood of  $\boldsymbol{\eta}$ . If  $E\|\underline{\mathbb{T}}(\mathbf{X})\|^4 < \infty$  then

$$E[\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}] = \frac{1}{n} \mathbf{a}(\boldsymbol{\theta}) + \frac{1}{n^2} \mathbf{b}(\boldsymbol{\theta}) + O\left(\frac{1}{n^3}\right), \quad (2.2.23)$$

where  $\mathbf{a}(\boldsymbol{\theta})$  and  $\mathbf{b}(\boldsymbol{\theta})$  are unknown functions of  $\boldsymbol{\theta}$ .

**Proof.**

Let  $\boldsymbol{\eta}$  and  $\mathbf{H}(\boldsymbol{\eta})$  be as defined in proposition 1, and define the sample mean of  $\underline{\mathbb{T}}(\mathbf{X})$  as

$$\bar{\underline{\mathbb{T}}}_n(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n \underline{\mathbb{T}}(\mathbf{X}_i).$$

To obtain the central moment of  $\bar{\underline{\mathbb{T}}}_n = \bar{\underline{\mathbb{T}}}_n(\mathbf{X})$ , let us first define the corresponding central moments of the elements  $\bar{\mathbb{T}}_{nj}(\mathbf{X})$  of  $\bar{\underline{\mathbb{T}}}_n$ , as the product moment of the variates about their respective means  $\eta_j$ 's and note that the expected value of the a random vector matrix is the expected value of it's random elements.

Writing  $\bar{\underline{\mathbb{T}}}_n$  for  $\bar{\underline{\mathbb{T}}}_n(\mathbf{X}_i)$  the first central moment of  $\bar{\mathbb{T}}_{nj}$  is,

$$E[\bar{\mathbb{T}}_{nj} - \eta_j] = \frac{1}{n} \sum_{i=1}^n E[\underline{\mathbb{T}}_j(\mathbf{X}_i) - \eta_j] = 0 \text{ so that}$$

$$E[\bar{\underline{\mathbb{T}}}_n - \boldsymbol{\eta}] = \mathbf{0}. \quad (2.2.24)$$

The second central moment of  $\bar{\mathbb{T}}_{nj}$  is

$$\mathbf{E}[\bar{\mathbf{T}}_{nj} - \eta_j][\bar{\mathbf{T}}_{nj'} - \eta_{j'}] = \frac{1}{n^2} \sum_{i=1}^n \mathbf{E}[\mathbf{T}_j(\mathbf{X}_i) - \eta_j][\mathbf{T}_{j'}(\mathbf{X}_i) - \eta_{j'}] = \frac{1}{n} \sigma_{jj'},$$

so that

$$\mathbf{E}[\bar{\mathbf{T}}_n - \boldsymbol{\eta}][\bar{\mathbf{T}}_n - \boldsymbol{\eta}]^T = \frac{1}{n} \boldsymbol{\Sigma}_T, \quad (2.2.25)$$

where  $\boldsymbol{\Sigma}_T = (\sigma_{jj'})$   $j, j' = 1, 2, \dots, k$  is the covariance matrix of  $\bar{\mathbf{T}}(\mathbf{X})$ .

For the third central moment the  $\bar{\mathbf{T}}_n$  we have

$$\begin{aligned} & \mathbf{E}[\bar{\mathbf{T}}_{nj} - \eta_j][\bar{\mathbf{T}}_{nj'} - \eta_{j'}][\bar{\mathbf{T}}_{nj''} - \eta_{j''}] \\ &= \frac{1}{n^3} \mathbf{E}\left\{ \sum_{i=1}^n [\mathbf{T}_j(\mathbf{X}_i) - \eta_j][\mathbf{T}_{j'}(\mathbf{X}_i) - \eta_{j'}][\mathbf{T}_{j''}(\mathbf{X}_i) - \eta_{j''}] \right\} \\ &= \frac{1}{n^2} \mathbf{E}[\mathbf{T}_j(\mathbf{X}_i) - \eta_j][\mathbf{T}_{j'}(\mathbf{X}_i) - \eta_{j'}][\mathbf{T}_{j''}(\mathbf{X}_i) - \eta_{j''}] \\ &= \frac{1}{n^2} \gamma_{jj'j''} \end{aligned}$$

Note that the  $\mathbf{X}_i$ 's are independent and hence  $\mathbf{E}[\mathbf{T}_j(\mathbf{X}_i) - \eta_j][\mathbf{T}_{j'}(\mathbf{X}_{i'}) - \eta_{j'}] = 0$ ,  $i \neq i'$  and  $\mathbf{E}[\mathbf{T}_j(\mathbf{X}_i) - \eta_j][\mathbf{T}_{j'}(\mathbf{X}_{i'}) - \eta_{j'}][\mathbf{T}_{j''}(\mathbf{X}_{i''}) - \eta_{j''}] = 0$ , for  $i \neq i' \neq i''$ .

Hence,

$$\mathbf{E}[\bar{\mathbf{T}}_n - \boldsymbol{\eta}][\bar{\mathbf{T}}_n - \boldsymbol{\eta}]^T[\bar{\mathbf{T}}_n - \boldsymbol{\eta}] = \frac{1}{n^2} \mathbf{J}_T, \quad (2.2.26)$$

where  $\underline{J}_T = (\gamma_{jj'j''})$   $j, j', j'' = 1, \dots, k$  is the third central moment of  $\underline{T}(X)$ .

Similarly the fourth central moment of  $\bar{\underline{T}}_n(X)$  is

$$E[(\bar{\underline{T}}_n - \underline{\eta})(\bar{\underline{T}}_n - \underline{\eta})^T(\bar{\underline{T}}_n - \underline{\eta})(\bar{\underline{T}}_n - \underline{\eta})^T] = \frac{1}{n^2} \underline{K}_T, \quad (2.2.27)$$

where  $\underline{K}_T = (\lambda_{jj'j''j'''})$   $j, j', j'', j''' = 1, \dots, k$  is the fourth central moment of  $\underline{T}(X)$

Similar results of the central moments of order  $k$  of the sample mean in a one dimensional variate (Cramer, 1946) are

$$E[(\bar{X}_n - E(\bar{X}_n))^{2k}] = \frac{1}{n^k} a_k(\theta) + O\left(\frac{1}{n^{k+1}}\right),$$

$$E[(\bar{X}_n - E(\bar{X}_n))^{2k-1}] = \frac{1}{n^k} b_k(\theta) + O\left(\frac{1}{n^{k+1}}\right) \quad (2.2.28)$$

Now consider the following Taylor's expansion of  $\mathbf{H}(\bar{\underline{T}}_n)$  about  $\underline{\eta} = E[\bar{\underline{T}}_n]$ :

$$\begin{aligned} \hat{\theta}_n - \theta &= \mathbf{H}(\bar{\underline{T}}_n) - \mathbf{H}(\underline{\eta}) \\ &= \frac{\partial \mathbf{H}(\underline{\eta})}{\partial \underline{\eta}} (\bar{\underline{T}}_n - \underline{\eta}) + \frac{1}{2} (\bar{\underline{T}}_n - \underline{\eta})^T \frac{\partial^2 \mathbf{H}(\underline{\eta})}{\partial \underline{\eta}^T \partial \underline{\eta}} (\bar{\underline{T}}_n - \underline{\eta}) \\ &\quad + \frac{1}{3!} (\bar{\underline{T}}_n - \underline{\eta})(\bar{\underline{T}}_n - \underline{\eta})^T \frac{\partial^3 \mathbf{H}(\underline{\eta})}{\partial \underline{\eta} \partial \underline{\eta}^T \partial \underline{\eta}} (\bar{\underline{T}}_n - \underline{\eta}) \end{aligned}$$

$$+ \frac{1}{4!} ((\bar{\mathbb{T}}_n - \boldsymbol{\eta})(\bar{\mathbb{T}}_n - \boldsymbol{\eta})^T)^2 \frac{\partial^4 \mathbf{H}(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}^T \partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T \partial \boldsymbol{\eta}} (\bar{\mathbb{T}}_n - \boldsymbol{\eta})(\bar{\mathbb{T}}_n - \boldsymbol{\eta})^T + \gamma_n, \quad (2.2.29)$$

where  $\gamma_n$  is a random variable which represents the remainder term. Let  $\mathbf{H}(\boldsymbol{\eta}) = (h_1(\boldsymbol{\eta}), \dots, h_l(\boldsymbol{\eta}), \dots, h_k(\boldsymbol{\eta}))^T$ . Consider Taylor's expansion of a particular element  $h_l(\bar{\mathbb{T}}_n)$  of  $\mathbf{H}(\bar{\mathbb{T}}_n)$ , about  $\boldsymbol{\eta}$  and, similar to (2.2.29), we have

$$\begin{aligned} & h_l(\bar{\mathbb{T}}_n) - h_l(\boldsymbol{\eta}) \\ &= \frac{\partial h_l(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} (\bar{\mathbb{T}}_n - \boldsymbol{\eta}) + \frac{1}{2} (\bar{\mathbb{T}}_n - \boldsymbol{\eta})^T \frac{\partial^2 h_l(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}^T \partial \boldsymbol{\eta}} (\bar{\mathbb{T}}_n - \boldsymbol{\eta}) \\ & \quad + \frac{1}{3!} (\bar{\mathbb{T}}_n - \boldsymbol{\eta})(\bar{\mathbb{T}}_n - \boldsymbol{\eta})^T \frac{\partial^3 h_l(\boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T \partial \boldsymbol{\eta}} (\bar{\mathbb{T}}_n - \boldsymbol{\eta}) \\ & \quad + \frac{1}{4!} (\bar{\mathbb{T}}_n - \boldsymbol{\eta})(\bar{\mathbb{T}}_n - \boldsymbol{\eta})^T \frac{\partial^4 h_l(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}^T \partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T \partial \boldsymbol{\eta}} (\bar{\mathbb{T}}_n - \boldsymbol{\eta})(\bar{\mathbb{T}}_n - \boldsymbol{\eta})^T + \gamma_{nl}. \end{aligned} \quad (2.2.30)$$

Note the expected value of the first term on the RHS of (2.2.30) is zero. To study the second and subsequent terms on the RHS of (2.2.30), we introduce the following notation:

Let

$$\mathbf{y} = \bar{\mathbb{T}}_n - \boldsymbol{\eta}, \quad \mathbf{B} = \mathbf{y}\mathbf{y}^T = (b_{ji}),$$

where  $b_{ji} = (\bar{\mathbb{T}}_{nj} - \eta_j)(\bar{\mathbb{T}}_{ni} - \eta_i)$  and

$$\mathbf{A} = \frac{\partial^2 h_l(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}^T \partial \boldsymbol{\eta}} = (a_{ij}^{(l)})_{i,j=1,\dots,k}$$

where  $a_{ji}^{(l)} = \frac{\partial^2 h_l(\boldsymbol{\eta})}{\partial \eta_j \partial \eta_i}$ .

We introduce the following quantities which will be used in evaluating the bias of  $\widehat{\boldsymbol{\theta}}_n$  from (2.2.29).

$$a_{rij}^{(l)} = \left( \frac{\partial a_{ij}^{(l)}}{\partial \boldsymbol{\eta}} \right)_r = \frac{\partial^3 h_l(\boldsymbol{\eta})}{\partial \eta_r \partial \eta_j \partial \eta_i},$$

$$a_{srij}^{(l)} = \left( \frac{\partial a_{rij}^{(l)}}{\partial \boldsymbol{\eta}^T} \right)_s = \frac{\partial^4 h_l(\boldsymbol{\eta})}{\partial \eta_s \partial \eta_r \partial \eta_j \partial \eta_i},$$

and

$$c_r = \bar{\mathbb{T}}_{nr} - \eta_r.$$

Then the second term on the RHS of (2.2.30), which is a scalar is written as

$$\begin{aligned} & \mathbb{E}[(\bar{\mathbb{T}}_n - \boldsymbol{\eta})^T \frac{\partial^2 h_l(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}^T \partial \boldsymbol{\eta}} (\bar{\mathbb{T}}_n - \boldsymbol{\eta})] \\ &= \mathbb{E}[\text{Tr}(\underline{\mathbf{A}} \underline{\mathbf{B}})] = \mathbb{E}[\sum_{i=1}^k \sum_{j=1}^k a_{ij}^{(l)} b_{ji}] \\ &= \frac{1}{n} [\sum_{i=1}^k \sum_{j=1}^k a_{ij}^{(l)} \sigma_{ji}]. \end{aligned} \tag{2.2.31}$$

The last equality is a consequence of (2.2.25), noting that the  $a_{ij}^{(l)}$ 's are fixed.

The subsequent higher order terms on the RHS of (2.2.30) are obtained by differentiating the second term with respect to  $\boldsymbol{\eta}$  and taking the expectation as follows. The expected value of the third term in (2.2.30) becomes

$$\mathbb{E}[(\bar{\mathbb{T}}_n - \boldsymbol{\eta})(\bar{\mathbb{T}}_n - \boldsymbol{\eta})^T \frac{\partial^3 h_l(\boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T \partial \boldsymbol{\eta}} (\bar{\mathbb{T}}_n - \boldsymbol{\eta})]$$

$$\begin{aligned}
&= \mathbb{E} \left[ \sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k c_r a_{rij}^{(l)} b_{ji} \right] \\
&= \frac{1}{n^2} \left[ \sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k a_{rij}^{(l)} \gamma_{rji} \right], \tag{2.2.32}
\end{aligned}$$

where the first equal sign in (2.2.32) is a result of taking the derivative of  $a_{ij}^{(l)}$  in (2.2.31), and  $\gamma_{rji}$  are given by (2.2.26).

From (2.2.26) we have  $\mathbb{E}[c_r a_{rij}^{(l)} b_{ji}] = O(n^{-2})$  and since  $k$  is fixed (2.2.32) is  $O(n^{-2})$ .

The expected value of the fourth term in (2.2.30) is given as,

$$\begin{aligned}
&\mathbb{E}[(\bar{\mathbb{T}}_n - \boldsymbol{\eta})(\bar{\mathbb{T}}_n - \boldsymbol{\eta})^T \frac{\partial^4 h_l(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}^T \partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T \partial \boldsymbol{\eta}} (\bar{\mathbb{T}}_n - \boldsymbol{\eta})(\bar{\mathbb{T}}_n - \boldsymbol{\eta})^T] \\
&= \mathbb{E} \left[ \sum_{s=1}^k \sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k c_s c_r a_{srij}^{(l)} b_{ji} \right] \\
&= \frac{1}{n^2} \left[ \sum_{s=1}^k \sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k a_{srij}^{(l)} \lambda_{srji} \right]. \tag{2.2.33}
\end{aligned}$$

From (2.2.27),  $\mathbb{E}[c_s c_r a_{srij}^{(l)} b_{ji}] = O(n^{-2})$  and similarly (2.2.33) is  $O(n^{-2})$ .

Note that the term with the fifth and sixth partial derivatives of  $h_l(\boldsymbol{\eta})$  will contribute  $O(n^{-3})$  to (2.2.30). So that from (2.2.30), as  $n \longrightarrow \infty$ , we obtain

$$\begin{aligned}
\mathbb{E}[h_l(\bar{\mathbb{T}}_n) - h_l(\boldsymbol{\eta})] &= \frac{\sum_{i=1}^k \sum_{j=1}^k a_{ij}^{(l)} \sigma_{ji}}{n} \\
&+ \frac{\sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k a_{rij}^{(l)} \gamma_{rji} + \sum_{s=1}^k \sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k a_{srij}^{(l)} \lambda_{srji}}{n^2} \\
&+ O(n^{-3}), \quad l = 1, \dots, k \tag{2.2.34}
\end{aligned}$$



and hence taking expectation of (2.2.29) we have

$$E[\hat{\theta}_n - \theta] = \frac{1}{n} \mathbf{a}(\theta) + \frac{1}{n^2} \mathbf{b}(\theta) + o\left(\frac{1}{n^2}\right), \quad \text{as } n \longrightarrow \infty, \quad (2.2.35)$$

where

$$\mathbf{a}(\theta) = \left( \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k a_{ij}^{(l)} \sigma_{ji} \right) \quad l = 1, \dots, k$$

and

$$\mathbf{b}(\theta) = \left( \frac{1}{3!} \sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k a_{rij}^{(l)} \gamma_{rji} + \frac{1}{4!} \sum_{s=1}^k \sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k a_{srji}^{(l)} \lambda_{srji} \right), \quad l = 1, \dots, k$$

are functions of  $\theta$  independent of  $n$ .

Note that  $\hat{\theta}_n$  will be unbiased for  $\theta$  if the  $k^{\text{th}}$  order differential of  $\mathbf{H}(\eta) \equiv \mathbf{0}$  for  $k > 1$ .

We show, in section 2.3.1, that the jackknifing method reduces the bias in the MLE by eliminating the leading term of order  $n^{-1}$  from the bias of the MLE (2.2.30). In case the leading term in the bias of  $\hat{\theta}_n$  is zero, the classical delete-one jackknife would not achieve the desired advantage of bias reduction in the estimator  $\hat{\theta}_n$ .

## 2.3 Jackknife Estimator of the MLE

Next we investigate the properties of the delete one observation in the jackknife estimation. The jackknife estimation was originally introduced by

Quenoulli as a method of bias reduction and was later extended by Tukey for robust interval estimation. Considerable research work has since been made to study the properties of the jackknife estimators in a variety of statistical estimation problems. Detailed studies of the technique appear in the literature of, among others, Schucany, Gray and Owen, 1971; Gray, Watkins and Adams, 1972; Miller, 1974; Sen, 1977.

The purpose of this chapter is to investigate the form of bias reduction achieved by the jackknife method applied to the maximum likelihood estimators of the parameters of an exponential family of distributions and furthermore, study the distribution of the jackknife estimators of the parameters in large samples. The results of this chapter are then applied to the parameter estimators of the length-biased version of the original exponential family of distributions, in the next chapter.

It should be noted that the jackknife method studied here is different from the one considered by Brillinger, 1964. Brillinger defines jackknife by dividing the sample observations into a fixed number of groups, and successively deleting the groups to obtain the pseudo-estimates. The asymptotic results are then obtained by increasing the sample size so that the number of observations in each group increases while the number of groups remains fixed. In our classical delete-one jackknife the large sample properties of the estimates the number of observations per group tends to infinity as the sample size increases without bound.

Given the maximum likelihood estimator  $\hat{\theta}_n$ , of  $\theta$  we define the pseudo-values

$$\tilde{\theta}_i = n\hat{\theta} - (n-1)\hat{\theta}^{(i)} \quad i = 1, \dots, n \quad (2.3.1)$$

where  $\hat{\theta}$  and  $\hat{\theta}^{(i)}$  are the MLE computed from all the  $n$  and  $n-1$  data points respectively, the superscript  $(i)$  denotes the  $i^{\text{th}}$  observation deleted from the data vector.

The classical Jackknife estimator of  $\theta$  defined as

$$\begin{aligned}\tilde{\theta}^J &= \frac{1}{n} \sum_{i=1}^n \tilde{\theta}_i \\ &= \hat{\theta} - \frac{(n-1)}{n} \sum_{i=1}^n (\hat{\theta}^{(i)} - \hat{\theta}).\end{aligned}\tag{2.3.2}$$

The Jackknife uses the estimates  $\hat{\theta}^{(i)}$  ( $i=1, \dots, n$ ) which are obtained by successively deleting one observation from the data and estimating the MLE from the remaining observations. So we need to study the behavior of the terms  $(\hat{\theta}^{(i)} - \hat{\theta}_n)$ , in order to establish the asymptotic results of the Jackknife estimators. The large sample properties of the Jackknife estimators are established through the use of estimating equations for the MLE.

Let  $U_{n-1}^{(i)}(\theta)$  be the solution of the estimating equation (2.2.3) with the  $i^{\text{th}}$  observation deleted ( $i=1, \dots, n$ ). Since  $\hat{\theta}$  is the MLE of  $\theta$  we have  $U_n(\hat{\theta}) = 0$  and similarly  $U_{n-1}^{(i)}(\hat{\theta}^{(i)}) = 0$ .

By definition

$$\begin{aligned}U_{n-1}^{(i)}(\theta) &= \sum_{\substack{j=1 \\ j \neq i}}^n \frac{\partial}{\partial \theta} \log f(X_j; \theta) \\ &= U_n(\theta) - \frac{\partial}{\partial \theta} \log f(X_i; \theta)\end{aligned}\tag{2.3.3}$$

so that

$$\begin{aligned}
U_{n-1}^{(i)}(\hat{\theta}) &= - \frac{\partial}{\partial \theta} \log f(X_i; \theta) \Big|_{\hat{\theta}} \\
&= - \mathbb{T}(X_i) - \psi(\hat{\theta}).
\end{aligned} \tag{2.3.4}$$

From (2.3.4) we write

$$\begin{aligned}
U_{n-1}^{(i)}(\hat{\theta}) - U_{n-1}^{(i)}(\hat{\theta}^{(i)}) &= - \frac{\partial}{\partial \theta} \log f(x_i; \theta) \Big|_{\hat{\theta}} + \frac{\partial}{\partial \theta} \log f(X_i; \theta) \Big|_{\hat{\theta}^{(i)}} \\
&= \psi(\hat{\theta}^{(i)}) - \psi(\hat{\theta})
\end{aligned} \tag{2.3.5}$$

and similarly using (2.2.6)

$$U_n(\hat{\theta}) - U_n(\hat{\theta}^{(i)}) = n[ \psi(\hat{\theta}) - \psi(\hat{\theta}^{(i)}) ]. \tag{2.3.6}$$

Since  $U_n(\hat{\theta}) = 0$ , we have

$$\mathbb{T}(X_i) + \psi(\hat{\theta}^{(i)}) = n[ \psi(\hat{\theta}^{(i)}) - \psi(\hat{\theta}) ] \tag{2.3.7}$$

which can be written as

$$\mathbb{T}(X_i) + \psi(\hat{\theta}) = (n-1)[ \psi(\hat{\theta}^{(i)}) - \psi(\hat{\theta}) ]. \tag{2.3.8}$$

Note that  $\mathbb{T}(X_i)$   $i = 1, \dots, n$  are iid with finite second moment, as a result

$$\max_{1 \leq k \leq n} \| \mathbb{T}(X_k) \| = o_p(\sqrt{n}), \tag{2.3.9}$$

hence

$$\max_{1 \leq i \leq n} \|\psi(\hat{\theta}^{(i)}) - \psi(\hat{\theta})\| = o_p\left(\frac{1}{\sqrt{n}}\right). \quad (2.3.10)$$

Given the strong consistency of  $\hat{\theta}$  and the continuity assumptions on  $\psi(\theta)$  in (2.2.13), we conclude from (2.3.10) that

$$\max_{1 \leq i \leq n} \|\hat{\theta}^{(i)} - \hat{\theta}\| = o_p\left(\frac{1}{\sqrt{n}}\right). \quad (2.3.11)$$

The following lemma will be used in the study of the asymptotic properties of the Jackknife estimate.

**Lemma 2.1.** Let  $\underline{h}_{n_i} = h(\hat{\theta}^{(i)} - \hat{\theta})$ , where  $0 < h < 1$ .

Then  $\max_{1 \leq i \leq n} \|\varphi(\hat{\theta} + \underline{h}_{n_i}) - \varphi(\hat{\theta})\| = o_p\left(\frac{1}{\sqrt{n}}\right)$ , uniformly in  $i$ .

**Proof.**

By (2.2.11)  $\max_{1 \leq i \leq n} \|\hat{\theta}^{(i)} - \hat{\theta}\| = o_p\left(\frac{1}{\sqrt{n}}\right)$ , so that uniformly in  $i$  ( $1 \leq i \leq n$ )

$$\varphi(\hat{\theta} + \underline{h}_{n_i}) = \varphi\left(\hat{\theta} + o_p\left(\frac{1}{\sqrt{n}}\right)\right) = \varphi(\hat{\theta}) + o_p\left(\frac{1}{\sqrt{n}}\right). \quad (2.3.12)$$

The last equality is due to the assumed continuity of  $\varphi(\theta)$  around  $\theta$  and almost sure convergence of  $\hat{\theta}$  to  $\theta$ .

□

**Lemma 2.2.** Let  $\hat{\theta}$  and  $\hat{\theta}^{(i)}$  be the MLE of  $\theta$  obtained from a complete sample and a delete one sample respectively. Then

$$\sum_{i=1}^n (\hat{\theta}^{(i)} - \hat{\theta}) = o_p\left(\frac{1}{n}\right). \quad (2.3.13)$$

**Proof.**

Consider Taylor's expansion of  $\psi(\hat{\theta}^{(i)})$  around  $\hat{\theta}$  and using (2.3.8) we have

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n [\mathbb{T}(X_i) + \psi(\hat{\theta})][\mathbb{T}(X_i) + \psi(\hat{\theta})]^T \\ &= \frac{(n-1)^2}{n} \sum_{i=1}^n [\psi(\hat{\theta}^{(i)}) - \psi(\hat{\theta})][\psi(\hat{\theta}^{(i)}) - \psi(\hat{\theta})]^T \\ &= \frac{(n-1)^2}{n} \varphi(\hat{\theta}^*) \sum_{i=1}^n (\hat{\theta}^{(i)} - \hat{\theta})(\hat{\theta}^{(i)} - \hat{\theta})^T \varphi(\hat{\theta}). \end{aligned} \quad (2.3.14)$$

Since the LHS of (2.3.14) is  $O_p(1)$  and the continuity assumption on  $\varphi(\theta)$ , (2.3.14) implies that,

$$\sum_{i=1}^n (\hat{\theta}^{(i)} - \hat{\theta})(\hat{\theta}^{(i)} - \hat{\theta})^T = O_p\left(\frac{1}{n}\right) \text{ and hence}$$

$$\sum_{i=1}^n \|\hat{\theta}^{(i)} - \hat{\theta}\| = O_p\left(\frac{1}{\sqrt{n}}\right). \quad (2.3.15)$$

Also by Taylor's expansion of  $\psi(\hat{\theta}^{(i)})$  around  $\hat{\theta}$  we write

$$\begin{aligned} \mathbf{0} &= \sum_{i=1}^n [\psi(\hat{\theta}^{(i)}) - \psi(\hat{\theta})] = \sum_{i=1}^n (\hat{\theta}^{(i)} - \hat{\theta})^T \varphi(\hat{\theta} + \underline{h}_{n_i}) \\ &= \sum_{i=1}^n (\hat{\theta}^{(i)} - \hat{\theta})^T \varphi(\hat{\theta}) + \sum_{i=1}^n (\hat{\theta}^{(i)} - \hat{\theta})^T [\varphi(\hat{\theta} + \underline{h}_{n_i}) - \varphi(\hat{\theta})] \end{aligned}$$

so that,

$$\begin{aligned}
\| \sum_{i=1}^n (\hat{\theta}_{-i} - \hat{\theta})^T \varphi(\hat{\theta}) \| &\leq \sum_{i=1}^n \| \hat{\theta}_{-i} - \hat{\theta} \|^T \max_i \| \varphi(\hat{\theta} + \frac{1}{n} \hat{\theta}_{-i}) - \varphi(\hat{\theta}) \| \\
&= O_p(\frac{1}{\sqrt{n}}) o_p(\frac{1}{\sqrt{n}}) = o_p(\frac{1}{n}),
\end{aligned} \tag{2.3.16}$$

which implies

$$\sum_{i=1}^n (\hat{\theta}^{(i)} - \hat{\theta}) = o_p(\frac{1}{n}).$$

□

### 2.3.1 Asymptotic Distribution and Bias of the Jackknife Estimator

Theorem 2.1 establishes that the Jackknife estimate  $\hat{\theta}^J$  and the MLE  $\hat{\theta}_n$  are equivalent in the sense that they both have the same asymptotic normal distribution.

**Theorem 2.1.** For  $\tilde{\theta}^J$  defined in (2.2.1)

$$\sqrt{n}(\tilde{\theta}^J - \theta) \longrightarrow_d N(0, I^{-1}(\theta)) \quad \text{as } n \rightarrow \infty \tag{2.3.17}$$

**Proof.**

From (2.3.1)

$$\tilde{\theta}^J - \hat{\theta} = - \frac{(n-1)}{n} \sum_{i=1}^n (\hat{\theta}^{(i)} - \hat{\theta}) \tag{2.3.18}$$

Since the MLE  $\hat{\theta}$  is asymptotically normal,  $\sqrt{n}(\hat{\theta} - \theta) \longrightarrow_d N(0, \Gamma^{-1}(\theta))$ , and by

Lemma 2.2, the RHS of (2.3.18) is  $o_p(\frac{1}{\sqrt{n}})$ , which implies  $\sqrt{n}(\tilde{\theta}^J - \hat{\theta}) = o_p(\frac{1}{\sqrt{n}})$ , and

using Slutskys Theorem we conclude  $\sqrt{n}(\tilde{\theta}^J - \theta) \longrightarrow_d N(0, \Gamma^{-1}(\theta))$ .

□

Since the Jackknife estimator  $\tilde{\theta}^J$  is a linear combination of the MLE  $\hat{\theta}_n$  from (2.3.1) the bias of the Jackknife estimator is expressed as:

$$\begin{aligned} E[\tilde{\theta}^J - \theta] &= E[\hat{\theta}_n - \theta] - \frac{(n-1)}{n} \sum_{i=1}^n E[(\hat{\theta}^{(i)} - \hat{\theta}_n)] \\ &= nE[\hat{\theta}_n - \theta] - (n-1)E[(\hat{\theta}_{n-1} - \theta)] \end{aligned} \quad (2.3.19)$$

From (2.2.26)  $nE[\hat{\theta}_n - \theta] = a(\theta) + \frac{1}{n} b(\theta) + O(\frac{1}{n^2})$  so that (2.3.17) becomes

$$E[\tilde{\theta}^J - \theta] = -\frac{1}{n(n-1)} b(\theta) + o(n^{-2}) \quad (2.3.20)$$

Comparing the bias of the MLE  $\hat{\theta}_n$  in (2.2.24) with the bias of the Jackknife estimator  $\tilde{\theta}^J$  in (2.3.20) we note that the Jackknife estimate has reduced the bias term of order  $\frac{1}{n}$ . This is exactly the purpose for which the Jackknife was proposed by Quenouille. We also noted that if  $a(\theta)$  vanishes to zero then the Jackknife estimator will not reduce the bias.

Gray, Watkins and Adams, 1972; Schucany, Gray and Owen, 1971 proposed



higher order Jackknife where, instead of delete-one,  $k$  observations are deleted successively and the statistic is recomputed from the remaining observations. For example if instead of deleting one observation, two observations are deleted, then the additional bias term of order  $n^{-2}$  will be eliminated.

Theorem 2.1 establishes that the variance of the Jackknife estimator  $\tilde{\theta}^J$  is a consistently estimates the variance of  $\hat{\theta}_n$ . Tukey proposed the variance of the sample variance of the pseudo values (2.2.1) divided by  $n$  would consistently estimate the variance of  $\tilde{\theta}^J$ . This is shown in Theorem 2.2.

**Theorem 2.2.** Consider the Jackknifed dispersion matrix

$$S_{nJ} = \frac{1}{n-1} \sum_{i=1}^n (\tilde{\theta}_i - \tilde{\theta}^J)(\tilde{\theta}_i - \tilde{\theta}^J)^T \quad (2.3.21)$$

Then  $S_{nJ} \longrightarrow I^{-1}(\theta)$ , as  $n \longrightarrow \infty$

where  $I(\theta)$  is the Fisher information matrix defined by A2.

**Proof.**

Consider the variates  $(\tilde{\theta}_i - \tilde{\theta}^J)$ , and using (2.3.10) and (2.3.11) substitute the expression (2.3.10) for  $\tilde{\theta}_i$ , which gives

$$(\tilde{\theta}_i - \tilde{\theta}^J) = (n-1) \left\{ (\hat{\theta}^{(i)} - \hat{\theta}) - \frac{1}{n} \sum_{j=1}^n (\hat{\theta}^{(j)} - \hat{\theta}) \right\}, \quad (2.3.22)$$

so that

$$\begin{aligned} S_{Jn} &= \frac{1}{n-1} \sum_{i=1}^n (\tilde{\theta}_i - \tilde{\theta}^J) (\tilde{\theta}_i - \tilde{\theta}^J)^T \\ &= (n-1) \sum_{i=1}^n (\hat{\theta}^{(i)} - \hat{\theta})(\hat{\theta}^{(i)} - \hat{\theta})^T - n(n-1) \left\{ \frac{1}{n} \sum_{i=1}^n (\hat{\theta}^{(i)} - \hat{\theta}) \right\} \left\{ \frac{1}{n} \sum_{j=1}^n (\hat{\theta}^{(j)} - \hat{\theta}) \right\}^T \end{aligned}$$

(2.3.23)

From Lemma 2.2,  $\sum_{i=1}^n (\hat{\theta}^{(i)} - \hat{\theta}) = o_p\left(\frac{1}{\sqrt{n}}\right)$ , so that  $n$  times the second term on the

RHS of (2.3.23) is  $o_p(1)$ .

To study the asymptotic behavior of the first term on RHS of (2.3.22) we use Taylor's expansion of  $\psi(\hat{\theta}^{(i)})$  around  $\hat{\theta}$ , together with Lemma 2.1 and write

$$\psi(\hat{\theta}^{(i)}) - \psi(\hat{\theta}) = \varphi(\hat{\theta})(\hat{\theta}^{(i)} - \hat{\theta}) + o_p(n^{-1/2}) \quad (2.3.24)$$

so that

$$\begin{aligned} & (n-1) \sum_{i=1}^n [\psi(\hat{\theta}^{(i)}) - \psi(\hat{\theta})][\psi(\hat{\theta}^{(i)}) - \psi(\hat{\theta})]^T \\ &= (n-1) \left\{ \varphi(\hat{\theta}) \left( \sum_{i=1}^n (\hat{\theta}^{(i)} - \hat{\theta})(\hat{\theta}^{(i)} - \hat{\theta})^T \right) \varphi(\hat{\theta}) + o_p(n^{-1}) \right\} \\ &= \varphi(\hat{\theta}) S_{nJ} \varphi(\hat{\theta}) + o_p(1). \end{aligned} \quad (2.3.25)$$

Now  $\varphi(\hat{\theta}) \rightarrow \varphi(\theta)$  a.s. and from (2.3.6) the LHS of (2.3.25) is given by

$$\frac{1}{n-1} \sum_{i=1}^n [\underline{\mathbb{T}}(X_i) + \psi(\theta)] [\underline{\mathbb{T}}(X_i) + \psi(\theta)]^T. \quad (2.3.26)$$

Since by Khintchine (WLLN) the term in (2.3.26) converges to  $I(\theta) = -\varphi(\theta)$ ,

then,

$$S_{nJ} \longrightarrow -[\varphi(\theta)]^{-1} \varphi(\theta) [\varphi(\theta)]^{-1} = I^{-1}(\theta), \quad \text{as } n \rightarrow \infty.$$

So that the sample variance of the pseudo values is a consistent estimator of the Jackknife variance estimator.

## Chapter 3.

# Parametric Regression Model: Exponential Family

### 3.1 Introduction.

This chapter provides a study of the Jackknife estimator of the regression parameters in a multiple regression setup. The bias of the Jackknife estimator of the regression parameters is derived and compared with the bias of the corresponding MLE. It is shown the Jackknife estimator and the MLE of the regression coefficients have the same asymptotic normal distribution. The convergence of the Jackknife variance estimate to the variance of the MLE is also studied.

We consider the statistical estimation in a linear regression model, where the mean of the response variate is a linear function of some fixed and known set of covariates. The notations and assumptions of the model are given below:

Let  $\mathbf{X} = (X_1, X_2, \dots, X_n)^T$  denote an  $n \times 1$  vector of response r.v.'s

$\boldsymbol{\mu} = (\mu_1, \mu_2, \dots, \mu_n)^T$  be the corresponding vector of expected value parameters.

$\mathbf{z}_i = (z_{i1}, z_{i2}, \dots, z_{ik})^T$  a vector of nonstochastic covariates,

$\mathbf{Z} = (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n)^T$  denote a  $n \times k$  design matrix with full rank  $k < n$

$\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_k)^T$  is  $k \times 1$  vector of unknown parameters to be estimated.

Assume  $(X_i, \mathbf{z}_i)$ ,  $i = 1, \dots, n$ , are the observations on the  $i^{\text{th}}$  subject. In some

applications it may be of interest to express the mean of  $X_i$  or some function of  $X_i$ , say  $T(X_i)$ , as

$$E[T(X_i)] = \mu_i = \mu(z_i^T \beta, \phi), \quad i = 1, \dots, n \quad (3.1)$$

The usual method of estimation in the univariate linear models is the least squares method. Least squares estimators are generally easier to handle mathematically and computationally, however they become less efficient when the underlying distribution of the error terms is significantly different from normal. Other types of robust estimators for linear models appearing in the literature are: 1) the maximum likelihood type of estimators (M-estimators) and 2) estimators derived from ranks (R-estimators). M-estimators introduced by Huber, 1965 is defined as an estimator  $T_n$  for  $\beta$  to be the solution to the equations

$$\sum_{i=1}^n \rho(X_i - z_i^T T_n) = \min_{\theta \in \Theta} \sum_{i=1}^n \rho(x_i - z_i^T T_n)$$

where  $\Theta$  is in  $\mathbb{R}^k$ . Huber, 1964 advocated replacing  $\rho(y)$  by  $\psi(y) = \partial \rho / \partial y$ .

Particular forms of  $\rho(y)$  are

(i)  $\rho(y) = y^2$

(ii)  $\rho(y) = |y|$

(iii)  $\rho(y) = -\log f(y)$ , where  $f(y)$  is the density of  $Y$ .

Our discussion will be focused on the general linear model with the score function of  $\rho(y) = -\log(f(y))$ . Specifically the  $X_i$   $i = 1, 2, \dots, n$  are assumed to follow an exponential family of distributions given by

$$f(X_i; \beta, \phi) = c(X_i, \phi) \exp\left\{ \mathbf{b}(z_i^T \beta, \phi)^T T(X_i) + \Psi(\mathbf{b}(z_i^T \beta, \phi)) \right\}, \quad (3.2)$$

where

$$\mathbf{b}(z_i^T \beta, \phi) = \left( b_l(z_i^T \beta, \phi) \quad l = 1, 2, \dots, q \right)^T \text{ is a } q\text{-vector of independent quantities,}$$

and  $\phi$  is the scale parameter.

$\Psi(\cdot)$  and  $\mathbf{b}(\cdot)$  are known function and

$\mathbf{T}(X_i) = (T_1(X_i), \dots, T_q(X_i))^T$  is a  $q$ -vector of sufficient statistics.

The specification of (3.2) includes a variety of regression models such as the log-linear models, survival data models and linear models, (see McCullagh and Nelder, (1989). Statistical analysis of linear regression models often involves estimation of the parameters  $\beta$  and making statistical inference about the parameter. We focus on estimating the parameters  $\beta$  via Jackknife estimation of the MLE  $\hat{\beta}$ . Both the bias and asymptotic distribution of  $\hat{\beta}_n$  and  $\hat{\beta}^J$  (the Jackknife estimator) are studied.

To derive the maximum likelihood estimators of  $\beta$  and  $\phi$ , we impose the same set of assumptions, "A" in Chapter 2, on the density  $f(X; \beta, \phi)$  in (3.2) with  $(\beta, \phi)$  replacing  $\theta$ , in order to derive asymptotic properties of the MLE  $\hat{\beta}_n$  of  $\beta$ . For simplicity of notation we write  $\beta^0$  for  $(\beta, \phi)$  and  $\mathbf{b}(z_i, \beta^0)$  for  $\mathbf{b}(z_i^T \beta, \phi)$ .

### Assumptions

(i). for  $l, j = 1, \dots, (k+1)$

$\frac{\partial}{\partial \beta_l^0} \log f(X_i; \beta^0)$  and  $\frac{\partial^2}{\partial \beta_l^0 \partial \beta_j^0} \log f(X_i; \beta^0)$  exist and are bounded by integrable

functions,

i.e.  $|\frac{\partial}{\partial \beta_l^0} \log f(X_i; \beta^0)| \leq K_1(x)dx$  and  $|\frac{\partial^2}{\partial \beta_l^0 \partial \beta_j^0} \log f(X_i; \beta^0)| \leq K_2(x)$

where  $\int K_i(x)dx < \infty \quad i = 1, 2$

(ii). The Fisher Information

$I_i(\beta^0) = E_{\beta^0} \left[ \frac{\partial}{\partial \beta^0} \log f(X_i; \beta^0) \right] \left[ \frac{\partial}{\partial \beta^0} \log f(X_i; \beta^0) \right]^T$  exists and is positive definite.

(iii).  $-E \left[ \frac{\partial^2}{\partial \beta^{0T} \partial \beta^0} \log f(X_i; \beta^0) \right] = I_i(\beta^0)$  is a positive definite matrix.

In addition we assume

4.  $E \left[ \|\underline{T}(X_i)\|^2 \right] < \infty$ , for all  $i = 1, \dots, n$  and all  $\beta^0 \in \Theta$

5.  $\mathbf{b}(z_i; \beta^0)$  is a  $q$ -vector of linearly independent components that are continuously differentiable, with continuous partial derivatives of  $\mathbf{b}(z_i; \beta^0)$  in some neighborhood of  $\beta^0$ . Also assume

$\left\| \frac{\partial}{\partial \beta, \phi} \mathbf{b}(z_i; \beta^0) \right\|$  is nonsingular

6.  $\Psi(\mathbf{b}(z_i; \beta^0))$  is continuously differentiable with respect to  $\beta^0$  and all order derivatives of  $\Psi(\cdot)$  are continuous in the neighborhood of  $\beta^0$ .

Note that the moments of  $\underline{T}(X_i)$  can be obtained from (3.2), we derive the first moment and the second moment of  $\underline{T}(X_i)$  by differentiating  $\int f(X_i; \beta^0) dx$  with respect to  $\beta^0$ .

$$0 = \int \frac{\partial}{\partial \beta} f(X_i; \beta^0) dx = \int (\mathbf{b}'(z_i; \beta^0))^T [\underline{T}(X_i) + \psi(\mathbf{b}(z_i; \beta^0))] f(X_i; \beta^0) dx$$

where

$$\mathbf{b}'(z_i; \beta^0)^T = \frac{\partial \mathbf{b}(z_i; \beta^0)^T}{\partial \beta^0} \quad \text{is a } (k+1) \times q \text{ vector,}$$

$$\psi(\mathbf{b}(z_i; \beta^0)) = \frac{\partial \Psi(\mathbf{b}(z_i; \beta^0))}{\partial \mathbf{b}(z_i; \beta^0)} \quad \text{is a } (q \times 1) \text{ vector}$$

so that,

$$E[\underline{T}(X_i)] = -\psi(\mathbf{b}(z_i, \beta^0)) \quad i = 1, \dots, n. \quad (3.3)$$

By differentiating  $\int \frac{\partial}{\partial \beta^0} f(X_i; \beta^0) dx$  again with respect to  $\beta^{0T}$  we obtain,

$$\begin{aligned} \mathbf{0} &= \int \frac{\partial^2}{\partial \beta^0 \partial \beta^{0T}} f(X_i; \beta^0) dx \\ &= \int \{(\mathbf{b}'(z_i, \beta^0))^T [\underline{T}(X_i) + \psi(\mathbf{b}(z_i, \beta^0))][\underline{T}(X_i) + \psi(\mathbf{b}(z_i, \beta^0))]^T \mathbf{b}'(z_i, \beta^0) f(X_i; \beta^0) dx \\ &\quad + \int [(\mathbf{b}''(z_i, \beta^0))^T [\underline{T}(X_i) + \psi(\mathbf{b}(z_i, \beta^0))] f(X_i; \beta^0) dx \\ &\quad + (\mathbf{b}'(z_i, \beta^0))^T \varphi(\mathbf{b}(z_i, \beta^0)) \mathbf{b}'(z_i, \beta^0), \end{aligned} \quad (3.4)$$

where

$$\varphi(\mathbf{b}(z_i, \beta^0)) = \frac{\partial^2 \Psi(\mathbf{b}(z_i, \beta^0))}{\partial (\mathbf{b}(z_i, \beta^0))^T \partial \mathbf{b}(z_i, \beta^0)}.$$

It follows from (3.4) and (3.3) that

$$\begin{aligned} E[(\mathbf{b}'(z_i, \beta^0))^T [\underline{T}(X_i) + \psi(\mathbf{b}(z_i, \beta^0))][\underline{T}(X_i) + \psi(\mathbf{b}(z_i, \beta^0))]^T \mathbf{b}'(z_i, \beta^0) \\ = -(\mathbf{b}'(z_i, \beta^0))^T \varphi(\mathbf{b}(z_i, \beta^0)) \mathbf{b}'(z_i, \beta^0), \end{aligned} \quad (3.5)$$

hence the covariance matrix of  $\underline{T}(X_i)$  is given by the  $(q \times q)$  matrix

$$E[\underline{T}(X_i) - E\underline{T}(X_i)] [\underline{T}(X_i) - E\underline{T}(X_i)]^T$$

$$= -\varphi(\mathbf{b}(z_i, \beta^0)). \quad (3.6)$$

Higher order central moments of  $\underline{T}(X_i)$  can also be obtained from the (3.2). The third and fourth central moments of  $\underline{T}(X_i)$  are given by

$$\begin{aligned} E[\underline{T}(X_i) - E\underline{T}(X_i)]^T [\underline{T}(X_i) - E\underline{T}(X_i)] [\underline{T}(X_i) - E\underline{T}(X_i)]^T \\ = -\Gamma(\mathbf{b}(z_i, \bar{\beta})) \end{aligned} \quad (3.7)$$

and

$$\begin{aligned} E[\underline{T}(X_i) - E\underline{T}(X_i)]^T [\underline{T}(X_i) - E\underline{T}(X_i)] [\underline{T}(X_i) - E\underline{T}(X_i)]^T [\underline{T}(X_i) - E\underline{T}(X_i)] \\ = 3 \varphi(\mathbf{b}(z_i, \beta^0)) \varphi(\mathbf{b}(z_i, \beta^0))^T - \Upsilon(\mathbf{b}(z_i, \beta^0)), \end{aligned} \quad (3.8)$$

where

$$\begin{aligned} \Gamma(\mathbf{b}(z_i, \beta^0)) &= \frac{\partial^3 \Psi(\mathbf{b}(z_i, \beta^0))}{\partial(\mathbf{b}(z_i, \beta^0))^T \partial \mathbf{b}(z_i, \beta^0) \partial(\mathbf{b}(z_i, \beta^0))^T} \text{ and} \\ \Upsilon(\mathbf{b}(z_i, \beta^0)) &= \frac{\partial^4 \Psi(\mathbf{b}(z_i, \beta^0))}{\partial(\mathbf{b}(z_i, \beta^0))^T \partial \mathbf{b}(z_i, \beta^0) \partial(\mathbf{b}(z_i, \beta^0))^T \partial \mathbf{b}(z_i, \beta^0)}. \end{aligned}$$

### 3.2 Estimation of Parameters of the Regression Model

The maximum likelihood estimator of  $\beta^0$ , denoted by  $\beta_n^0$ , is a function of the



sample data, of size  $n$ , which maximizes the log-likelihood function. To study the maximum likelihood estimators we present the partial derivatives of the log-likelihood with respect to the parameters  $\beta^0$ .

The first and second partial derivatives of the  $\log f(X_i; \beta^0)$  in (3.2) with respect to  $\beta^0$  are obtain as:

$$\frac{\partial \log f(X_i; \beta^0)}{\partial \beta^0} = (\mathbf{b}'(\mathbf{z}_i; \beta^0))^T \{ \underline{\mathbb{T}}(X_i) + \psi(\mathbf{b}(\mathbf{z}_i; \beta^0)) \} \quad (3.9)$$

and

$$\begin{aligned} \frac{\partial^2 \log f(X_i; \beta^0)}{\partial \beta^0 \partial \beta^{0T}} &= (\mathbf{b}'(\mathbf{z}_i; \beta^0))^T \varphi(\mathbf{b}(\mathbf{z}_i; \beta^0)) \mathbf{b}'(\mathbf{z}_i; \beta^0) \\ &+ [(\mathbf{b}''(\mathbf{z}_i; \beta^0)) \underline{\mathbb{T}}(X_i) + \psi(\mathbf{b}(\mathbf{z}_i; \beta^0))]. \end{aligned} \quad (3.10)$$

Then from (3.3) and (3.10) the Fisher information from the  $i^{th}$  data point is

$$\begin{aligned} \mathbf{I}_i(\beta^0) &= E\{ (\mathbf{b}'(\mathbf{z}_i; \beta^0))^T [\underline{\mathbb{T}}(X_i) + \psi(\mathbf{b}(\mathbf{z}_i; \beta^0))] [\underline{\mathbb{T}}(X_i) + \psi(\mathbf{b}(\mathbf{z}_i; \beta^0))]^T \mathbf{b}'(\mathbf{z}_i; \beta^0) \} \\ &= - (\mathbf{b}'(\mathbf{z}_i; \beta^0))^T \varphi(\mathbf{b}(\mathbf{z}_i; \beta^0)) \mathbf{b}'(\mathbf{z}_i; \beta^0). \end{aligned} \quad (3.11)$$

The estimating equations for  $\beta^0$  are developed in the same fashion as in Chapter 2, and are given by

$$U_n(\beta^0) = \frac{\partial}{\partial \beta^0} \log L_n(\beta^0; \mathbf{X}, \underline{\mathbb{Z}})$$

$$= \sum_{i=1}^n (\mathbf{b}'(\mathbf{z}_i; \boldsymbol{\beta}^0))^T [\mathbb{T}(X_i) + \boldsymbol{\psi}(\mathbf{b}(\mathbf{z}_i; \boldsymbol{\beta}^0))]. \quad (3.12)$$

Note that the likelihood estimating equation (3.12) is now a  $(k+1) \times 1$  vector which can be written as

$$\begin{aligned} U_n(\boldsymbol{\beta}^0) &= \frac{\partial}{\partial \boldsymbol{\beta}^0} \sum_{i=1}^n \log f(X_i; \boldsymbol{\beta}^0) \\ &= (U_{n1}(\boldsymbol{\beta}^0), U_{n2}(\boldsymbol{\beta}^0))^T, \end{aligned}$$

where the first  $(k \times 1)$  rows are defined by

$$U_{n1}(\boldsymbol{\beta}^0) = \sum_{i=1}^n \partial \log f(X_i; \boldsymbol{\beta}^0) / \partial \boldsymbol{\beta}^0$$

and the  $(k+1)^{th}$  element is

$$U_{n2}(\boldsymbol{\beta}^0) = \sum_{i=1}^n \partial \log f(X_i; \boldsymbol{\beta}^0) / \partial \phi.$$

So that the MLE  $\hat{\boldsymbol{\beta}}^0$  is the solution to the equation (3.12) set equal to zero, and satisfies the equation

$$U_n(\hat{\boldsymbol{\beta}}_n^0) = \mathbf{0}. \quad (3.13)$$

In general, explicit closed-form expressions for the MLE  $\hat{\boldsymbol{\beta}}_n^0$  may not exist. In such cases iterative methods like the **Newton-Raphson** techniques for finding the MLE can be used as described below. Consider the first order Taylor's series expansion of the estimating equation around some initial guess  $\boldsymbol{\beta}_n^{0(0)}$ :

$$\begin{aligned} \mathbf{0} &= U_n(\boldsymbol{\beta}^0) \Big|_{\widehat{\boldsymbol{\beta}}_n^0} \\ &= U_n(\boldsymbol{\beta}^0) \Big|_{\boldsymbol{\beta}_n^{0(0)}} + (\widehat{\boldsymbol{\beta}}^0 - \boldsymbol{\beta}_n^{0(0)}) U'_n(\boldsymbol{\beta}^0) \Big|_{\boldsymbol{\beta}_n^*}, \end{aligned}$$

where  $\boldsymbol{\beta}_n^*$  is on the line segment joining  $\widehat{\boldsymbol{\beta}}_n^0$ ,  $\boldsymbol{\beta}_n^{0(0)}$  and  $U'_n(\boldsymbol{\beta}^0) = \frac{\partial U_n(\boldsymbol{\beta}^0)}{\partial \boldsymbol{\beta}^{0T}}$ .

Then we have

$$\widehat{\boldsymbol{\beta}}_n^0 = \boldsymbol{\beta}_n^{0(0)} - \left\{ U_n(\boldsymbol{\beta}^0) \Big|_{\boldsymbol{\beta}_n^{0(0)}} \right\} \left\{ U'_n(\boldsymbol{\beta}^0) \Big|_{\boldsymbol{\beta}_n^*} \right\}^{-1}$$

and by choosing  $\boldsymbol{\beta}_n^{0(0)}$  in some neighborhood of  $\widehat{\boldsymbol{\beta}}_n^0$  we may use the iteration procedure to obtain a first step estimator:

$$\widehat{\boldsymbol{\beta}}_n^{0(1)} = \widehat{\boldsymbol{\beta}}_n^{0(0)} - \left\{ U_n(\boldsymbol{\beta}^0) \Big|_{\widehat{\boldsymbol{\beta}}_n^{0(0)}} \right\} \left\{ U'_n(\boldsymbol{\beta}^0) \Big|_{\widehat{\boldsymbol{\beta}}_n^{0(0)}} \right\}^{-1}.$$

The process is repeated by updating  $\widehat{\boldsymbol{\beta}}_n^{0(1)}$ . So that at the  $i^{\text{th}}$  iteration estimates obtained from the previous  $(i-1)^{\text{th}}$  are used to refine the current estimates until the process is stopped when  $\widehat{\boldsymbol{\beta}}_n^{0(i)} - \widehat{\boldsymbol{\beta}}_n^{0(i-1)}$  is sufficiently small.

To derive the variance covariance matrix of  $\widehat{\boldsymbol{\beta}}_n^0$ , which is the inverse of the Fisher information matrix, we proceed as follows. Observe that

$$\begin{aligned} &U_n(\boldsymbol{\beta}^0) U_n(\boldsymbol{\beta}^0)^T \\ &= \sum_{i=1}^n (\mathbf{b}'(\mathbf{z}_i, \boldsymbol{\beta}^0))^T [\mathbb{T}(X_i) + \boldsymbol{\psi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0))] [\mathbb{T}(X_i) + \boldsymbol{\psi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0))]^T \mathbf{b}'(\mathbf{z}_i, \boldsymbol{\beta}^0) \end{aligned}$$

$$+ \sum_{i \neq i'} (\mathbf{b}'(\mathbf{z}_i, \boldsymbol{\beta}^0))^T [\underline{\mathbb{T}}(X_i) + \boldsymbol{\psi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0))] [\underline{\mathbb{T}}(X_{i'}) + \boldsymbol{\psi}(\mathbf{b}(\mathbf{z}_{i'}, \boldsymbol{\beta}^0))]^T \mathbf{b}'(\mathbf{z}_{i'}, \boldsymbol{\beta}^0) \quad (3.14)$$

Since the  $X_i$ 's are independent, taking expectation of the vector (3.14), and noting from (3.3) that  $E[\underline{\mathbb{T}}(X_i) + \boldsymbol{\psi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0))] = \mathbf{0}$ , the second term  $\equiv \mathbf{0}$ , hence we get

$$\begin{aligned} & \frac{1}{n} E[U_n(\boldsymbol{\beta}^0) U_n(\boldsymbol{\beta}^0)^T] \\ &= \frac{1}{n} E \sum_{i=1}^n (\mathbf{b}'(\mathbf{z}_i, \boldsymbol{\beta}^0))^T [\underline{\mathbb{T}}(X_i) + \boldsymbol{\psi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0))] [\underline{\mathbb{T}}(X_i) + \boldsymbol{\psi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0))]^T \mathbf{b}'(\mathbf{z}_i, \boldsymbol{\beta}^0) \\ &= -\frac{1}{n} \sum_{i=1}^n (\mathbf{b}'(\mathbf{z}_i, \boldsymbol{\beta}^0))^T \boldsymbol{\varphi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0)) \mathbf{b}'(\mathbf{z}_i, \boldsymbol{\beta}^0) \\ &= \bar{\mathbb{I}}_n(\boldsymbol{\beta}^0). \end{aligned} \quad (3.15)$$

where  $\bar{\mathbb{I}}_n(\boldsymbol{\beta}^0)$  is the average Fisher information matrix for the  $n$  data points which from (3.11) is also given by

$$\bar{\mathbb{I}}_n(\boldsymbol{\beta}^0) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_i(\boldsymbol{\beta}^0)$$

The second partial derivatives with respect to  $\boldsymbol{\beta}^0$  of the log-likelihood function is given by

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\beta}^{0T}} U_n(\boldsymbol{\beta}^0) &= \sum_{i=1}^n \frac{\partial^2}{\partial \boldsymbol{\beta}^0 \partial \boldsymbol{\beta}^{0T}} \log f(X_i; \boldsymbol{\beta}^0) \\ &= \sum_{i=1}^n (\mathbf{b}'(\mathbf{z}_i, \boldsymbol{\beta}^0))^T \boldsymbol{\varphi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0)) \mathbf{b}'(\mathbf{z}_i, \boldsymbol{\beta}^0) + (\mathbf{b}''(\mathbf{z}_i, \boldsymbol{\beta}^0))^T [\underline{\mathbb{T}}(X_i) + \boldsymbol{\psi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0))] \end{aligned} \quad (3.16)$$

Taking the expected value of the expression in (3.16) the second term is zero and we have

$$\begin{aligned} -\frac{1}{n}E\left[\frac{\partial}{\partial\beta^T}U_n(\bar{\beta})\right] &= -\frac{1}{n}\sum_{i=1}^n(\mathbf{b}'(\mathbf{z}_i,\beta^0))^T\varphi(\mathbf{b}(\mathbf{z}_i,\beta^0))\mathbf{b}'(\mathbf{z}_i,\beta^0) \\ &= \bar{\mathbf{I}}_n(\beta^0). \end{aligned} \quad (3.17)$$

Since  $E[U_n(\beta^0)] = \mathbf{0}$ ,  $\bar{\mathbf{I}}_n(\beta^0)$  is also the covariance matrix of  $U_n(\beta^0)$ .

Hence from (3.16) we write

$$-\frac{\partial}{\partial\beta^{0T}}U_n(\beta^0) = \sum_{i=1}^n \mathbf{I}_i(\beta^0) - \sum_{i=1}^n \left\{ [\mathbf{b}''(\mathbf{z}_i,\beta^0)[\underline{\mathbb{T}}(X_i) + \psi(\mathbf{b}(\mathbf{z}_i,\beta^0))]\right\} \quad (3.18)$$

To study the behavior of the second term in (3.16) in large samples, we assume

$$-\mathbf{b}''(\mathbf{z}_i,\beta^0)\varphi(\mathbf{b}(\mathbf{z}_i,\beta^0))\mathbf{b}''(\mathbf{z}_i,\beta^0)^T = \gamma_i(\beta^0) > 0, \quad (i = 1, \dots, n)$$

and

$$\frac{1}{n}\sum_{i=1}^n \gamma_i(\beta^0) \longrightarrow \underline{\mathbb{T}}(\beta^0) \quad (\text{p.d.}) \quad \text{as } n \longrightarrow \infty \quad (3.19)$$

Let

$$\mathbf{Q}_n(\beta^0) = -\sum_{i=1}^n \left\{ (\mathbf{b}''(\mathbf{z}_i,\beta^0))^T [\underline{\mathbb{T}}(X_i) + \psi(\mathbf{b}(\mathbf{z}_i,\beta^0))]\right\}$$

and note that  $\mathbf{Q}_n(\beta^0)$  is a sum of independent random variates centered around their respective means and

$$E\left[\frac{1}{n}Q_n(\beta^0)\right] = 0$$

and by (3.6) and (3.19) as  $n \rightarrow \infty$

$$\text{Cov}\left[\frac{1}{n}Q_n(\beta^0)\right] = -\frac{1}{n^2} \sum_{i=1}^n (\mathbf{b}''(\mathbf{z}_i, \beta^0))^T \varphi(\mathbf{b}(\mathbf{z}_i, \beta^0)) \mathbf{b}''(\mathbf{z}_i, \beta^0) \rightarrow 0 \quad (3.20)$$

So by Chebyshev's inequality, the second term in (3.16) divided by  $n$  is

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{b}''(\mathbf{z}_i, \beta^0))^T \left\{ \mathbb{T}(X_i) + \psi(\mathbf{b}(\mathbf{z}_i, \beta^0)) \right\} \xrightarrow{p} 0 \quad \text{as } n \rightarrow \infty \quad (3.21)$$

Hence the LHS of (3.18) divided by  $n$  is the sample average of the  $I_i(\beta^0)$  plus the remainder term which is asymptotically negligible.

To study the asymptotic behavior of  $\hat{\beta}_n^0$  we make following assumptions parallel to the assumptions in Chapter I about the behavior of the second partial derivative of the log-likelihood in some neighborhood of  $\beta^0$ . Note that the loglikelihood is a function of the  $\mathbf{b}'(\mathbf{z}_i, \beta^0)$ , hence we need some assumptions about the continuity behavior of the  $\mathbf{b}'(\mathbf{z}_i, \beta^0)$  and  $\psi(\mathbf{b}(\mathbf{z}_i, \beta^0))$  in some small neighborhood of  $\beta^0$ . In particular we assume for

$$\|\beta^{0*} - \beta^0\| < \delta, \text{ and } \delta \rightarrow 0$$

$$1a. \quad \sup \frac{1}{n} \sum_{i=1}^n \left\| \left\{ \mathbf{b}'(\mathbf{z}_i, \beta^0) \Big|_{\beta^{0*}} - \mathbf{b}'(\mathbf{z}_i, \beta^0) \right\} \mathbb{T}(X_i) \right\| \rightarrow 0$$

$$1b. \quad \sup \frac{1}{n} \sum_{i=1}^n \left\| (\mathbf{b}'(\mathbf{z}_i, \beta^0))^T \psi(\mathbf{b}(\mathbf{z}_i, \beta^0)) \Big|_{\beta^{0*}} - (\mathbf{b}'(\mathbf{z}_i, \beta^0))^T \psi(\mathbf{b}(\mathbf{z}_i, \beta^0)) \right\| \rightarrow 0$$

$$2a. \sup \frac{1}{n} \sum_{i=1}^n \left\| \left\{ \mathbf{b}''(\mathbf{z}_i, \boldsymbol{\beta}^0) \right\} \Big|_{\boldsymbol{\beta}^{0*}} - \mathbf{b}''(\mathbf{z}_i, \boldsymbol{\beta}^0) \right\| \mathbb{T}(X_i) \parallel \longrightarrow 0$$

$$2b. \sup \frac{1}{n} \sum_{i=1}^n \left\| \mathbf{b}''(\mathbf{z}_i, \boldsymbol{\beta}^0) \boldsymbol{\psi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0)) \Big|_{\boldsymbol{\beta}^{0*}} - \mathbf{b}''(\mathbf{z}_i, \boldsymbol{\beta}^0) \boldsymbol{\psi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0)) \right\| \parallel \longrightarrow 0$$

2c.

$$\sup \frac{1}{n} \sum_{i=1}^n \left\| \left\{ \left( \mathbf{b}'(\mathbf{z}_i, \boldsymbol{\beta}^0) \right)^T \boldsymbol{\varphi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0)) \mathbf{b}'(\mathbf{z}_i, \boldsymbol{\beta}^0) \Big|_{\boldsymbol{\beta}^{0*}} - \left( \mathbf{b}'(\mathbf{z}_i, \boldsymbol{\beta}^0) \right)^T \boldsymbol{\varphi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0)) \mathbf{b}'(\mathbf{z}_i, \boldsymbol{\beta}^0) \right\} \right\|$$

$\longrightarrow 0$

and

$$\begin{aligned} & \mathbb{E}_{\boldsymbol{\beta}^0} \left\{ \sup \frac{1}{n} \left\| \sum_{i=1}^n \left[ \left( \mathbf{b}''(\mathbf{z}_i, \boldsymbol{\beta}^0) \right)^T \left[ \mathbb{T}(X_i) + \boldsymbol{\psi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0)) \right] \right] \Big|_{\boldsymbol{\beta}^{0*}} \right. \right. \\ & \quad \left. \left. - \frac{1}{n} \sum_{i=1}^n \left[ \left( \mathbf{b}''(\mathbf{z}_i, \boldsymbol{\beta}^0) \right)^T \left[ \mathbb{T}(X_i) + \boldsymbol{\psi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0)) \right] \right] \right\| \right\} \longrightarrow 0 \quad \text{as } \delta \longrightarrow 0 \quad (3.22) \end{aligned}$$

We also make the assumption that as  $n \rightarrow \infty$ ,

$$\frac{1}{n} \sum_{i=1}^n \mathbf{I}_i(\boldsymbol{\beta}^0) = \bar{\mathbf{I}}_n(\boldsymbol{\beta}^0) \longrightarrow \mathbf{I}(\boldsymbol{\beta}^0). \quad (3.23)$$

So that by (3.22) and (3.23)

$$-\frac{1}{n} \frac{\partial}{\partial \boldsymbol{\beta}^{0T}} U_n(\boldsymbol{\beta}^0) \Big|_{\boldsymbol{\beta}^{0*}} \longrightarrow^p \mathbf{I}(\boldsymbol{\beta}^0). \quad (3.24)$$

Note that  $\frac{1}{n} U_n(\boldsymbol{\beta}^0)$  is the sum of independent variates having finite second moment  $\bar{\mathbf{I}}_n(\boldsymbol{\beta}^0)$  moreover the variates are centered around their means. Since the second moment of the  $\mathbb{T}(X_i)$  and hence of the  $U_n(\boldsymbol{\beta}^0)$  exist and is finite, we can applying the Liapouov Theorem to claim

$$\frac{1}{\sqrt{n}} U_n(\beta^0) \longrightarrow \mathfrak{D} \mathcal{N}(\mathbf{0}, \mathbf{I}(\beta^0)). \quad (3.25)$$

Next we consider the first order Taylor's expansion of  $U_n(\hat{\beta}_n^0)$  around  $\beta^0$ .

$$\mathbf{0} = \frac{1}{n} U_n(\hat{\beta}_n^0) = \frac{1}{n} U_n(\beta^0) + \frac{1}{n} \frac{\partial}{\partial \beta^0} U_n(\beta^0) \Big|_{\beta^{0*}} (\hat{\beta}_n^0 - \beta^0). \quad (3.26)$$

Furthermore since  $\frac{1}{\sqrt{n}} U_n(\beta^0)$  is  $O_p(1)$ ,  $\frac{1}{n} \frac{\partial}{\partial \beta^0} U_n(\beta^0) \Big|_{\beta^{0*}} = O_p(1)$  and the LHS of

(3.26) is  $\mathbf{0}$  and also using (3.24) we conclude that

$$(\hat{\beta}_n^0 - \beta^0) = o_p(1). \quad (3.27)$$

From (3.26) we have

$$\sqrt{n}(\hat{\beta}_n^0 - \beta^0) = \frac{1}{\sqrt{n}} U_n(\beta^0) \mathbf{I}^{-1}(\beta^0) + o(1). \quad (3.28)$$

By (3.25) the first term in the RHS of (3.28) is asymptotically normal and hence by the Slutsky's theorem LHS of (3.28) is

$$\sqrt{n}(\hat{\beta}_n^0 - \beta^0) \longrightarrow \mathfrak{D} \mathcal{N}(\mathbf{0}, \mathbf{I}^{-1}(\beta^0)). \quad (3.29)$$

### Special Case

In the classical linear model we have  $\mathbf{b}(\mathbf{z}_i^T \beta, \phi) = \mathbf{z}_i^T \beta$ ,  $-\psi(\mathbf{b}(\mathbf{z}_i^T \beta)) = \mathbf{z}_i^T \beta$  and  $\underline{\mathbb{T}}(X_i) = \mathbb{T}(X_i)$ , in this case  $\phi$  is a nuisance parameter which does not influence the



estimating equations for  $\beta$ , hence without loss of generality  $\phi$  can be taken to equal unity, so that the maximum likelihood estimating equations (3.12), for  $\beta$  becomes

$$\sum_{i=1}^n z_i [\mathbb{T}(X_i) - z_i^T \beta] \quad (3.30)$$

which leads to the MLE  $\hat{\beta}_n$  given by

$$\hat{\beta}_n = \left\{ \sum_{i=1}^n z_i z_i^T \right\}^{-1} \sum_{i=1}^n z_i \mathbb{T}(X_i). \quad (3.31)$$

**Example 3.1** Let  $\log X_1, \dots, \log X_n$  be i.i.d. r.v.'s following a normal distribution with mean  $z_i^T \beta$  and common variance  $\phi \equiv \sigma^2$  ( $i = 1, \dots, n$ ). Then the  $X_i$  have a lognormal density given by

$$\begin{aligned} f(X_i; \beta, \sigma^2) &= \frac{1}{\sqrt{2\pi} \sigma X_i} e^{-\frac{1}{\sigma^2} (\log X_i - z_i^T \beta)^2} \\ &= c(X) e^{\underline{\mathbb{T}}(X_i)^T \mathbf{b}(z_i^T \beta, \sigma^2) + \Psi(\mathbf{b}(z_i^T \beta, \sigma^2))}, \quad X_i > 0 \quad (i = 1, \dots, n) \end{aligned}$$

where

$$\underline{\mathbb{T}}(X_i) = (\mathbb{T}_1(X_i), \mathbb{T}_2(X_i))^T = (\log X_i, (\log X_i)^2)^T,$$

$$\mathbf{b}(z_i^T \beta, \sigma^2) = (b_1(z_i^T \beta, \sigma^2), b_2(z_i^T \beta, \sigma^2)) = \left( \frac{z_i^T \beta}{\sigma^2}, -\frac{1}{\sigma^2} \right)^T$$

and

$$\Psi(\mathbf{b}(z_i^T \beta, \sigma^2)) = -\frac{1}{2} \left( \frac{b_1(z_i^T \beta, \sigma^2)^2}{b_2(z_i^T \beta, \sigma^2)} + \log -2b_2(z_i^T \beta, \sigma^2) \right).$$

The partial derivatives of  $b(\mathbf{z}_i^T \boldsymbol{\beta}, \sigma^2)$  are

$$\mathbf{b}'(\mathbf{z}_i^T \boldsymbol{\beta}, \sigma^2) = \left( b'_1(\mathbf{z}_i^T \boldsymbol{\beta}, \sigma^2), b'_2(\mathbf{z}_i^T \boldsymbol{\beta}, \sigma^2) \right) = \left( \frac{\mathbf{z}_i^T \boldsymbol{\beta}}{\sigma^2}, -\frac{1}{\sigma^2} \right)$$

$$= \begin{bmatrix} \frac{1}{\sigma^2} \mathbf{z}_i & 0 \\ -\frac{1}{\sigma^4} \mathbf{z}_i^T \boldsymbol{\beta} & \frac{1}{2\sigma^4} \end{bmatrix}$$

and

$$\psi(\mathbf{b}(\mathbf{z}_i^T \boldsymbol{\beta}, \sigma^2)) = \frac{1}{2} \begin{bmatrix} \frac{b_1(\mathbf{z}_i^T \boldsymbol{\beta}, \sigma^2)}{b_2(\mathbf{z}_i^T \boldsymbol{\beta}, \sigma^2)} \\ \frac{b_1(\mathbf{z}_i^T \boldsymbol{\beta}, \sigma^2)^2}{2b_2(\mathbf{z}_i^T \boldsymbol{\beta}, \sigma^2)^2} + \frac{1}{b_2(\mathbf{z}_i^T \boldsymbol{\beta}, \sigma^2)} \end{bmatrix},$$

so that the MLE of  $(\boldsymbol{\beta}, \sigma^2)$  is a solution to the  $(k+1) \times 1$  estimating equations

$$\mathbf{U}_n(\boldsymbol{\beta}, \sigma^2) = \begin{bmatrix} \frac{1}{\sigma^2} \sum (\log X_i - \mathbf{z}_i^T \boldsymbol{\beta}) \mathbf{z}_i \\ \frac{1}{2\sigma^4} \sum (\log X_i - \mathbf{z}_i^T \boldsymbol{\beta})^2 - \frac{n}{2\sigma^2} \end{bmatrix} = \mathbf{0}.$$

### 3.2.1 Bias of the Regression Parameter Estimates

We have already shown, in chapter II, that the MLE are generally not

unbiased and that the leading term in the bias of the MLE is of order  $n^{-1}$ , which was reduced by jackknifing the MLE. In this section we establish the bias of the regression parameters and demonstrate the extent to which the classical Jackknife estimator can reduce the bias of the MLE of the regression coefficients.

Since the  $\mathbf{z}_i$  are most often not the same, the  $E[\mathbb{T}(X_i)]$  which are functions of the  $\mathbf{z}_i$ , would generally be different and may not be identically distributed. Then the estimating equations are a sum of independent but not identically distributed random variables. The parameter vector  $\beta^0$  enters in the estimating functions as a product of the  $\mathbb{T}(X_i)$ , and further the estimating functions are an implicit function of the vector  $\beta^0$ , hence some iteration procedures are needed to estimate the parameters. Given that we can not write an expression of  $\hat{\beta}_n^0$  explicitly from the estimating equations, we rely on Taylor expansion of the estimating equations, to study the large-sample bias in  $\hat{\beta}_n^0$ .

**Proposition 3.1** The MLE  $\hat{\beta}_n^0$  has an asymptotic bias of order  $n^{-1}$ . For simplicity of presentation we write  $U_n$  for  $U_n(\beta^0)$ . The  $j^{\text{th}}$  component of  $U_n(\beta^0)$  will be denoted by  $U_{nj}$  ( $j = 1, 2, \dots, k+1$ ).

$$E[\hat{\beta}_n^0 - \beta^0] = \frac{\mathbf{a}(\beta^0)}{n} + \frac{\mathbf{b}(\beta^0)}{n^2} + O_p\left(\frac{1}{n^2}\right)$$

where  $\mathbf{a}(\cdot)$  and  $\mathbf{b}(\cdot)$  are known functions of  $\beta^0$  and are independent of  $n$ .

**Proof.**

Note that the  $\mathbf{z}_i$  and  $\mathbf{b}$  are fixed and known vectors.

Assume there is a function  $\mathbf{H} : \mathbb{R}^{k+1} \rightarrow \mathbb{R}^{k+1}$  given by  $\mathbf{H}(U_n(\beta^0)) = \beta^0$ , so that

$\mathbf{H}$  is a mapping from the estimating equations space into the parameter space.

Let  $\mathbf{H} = (h_1, \dots, h_{k+1})$  where the elements  $h_l$  ( $l = 1, \dots, k+1$ ) of  $\mathbf{H}$  have continuous fourth order partial derivatives with respect to the vector  $U_n(\beta^0)$ .

We also assume the fourth moment of  $U_n(\beta^0)$  exist and is finite, that is,

$$E\|U_n(\beta^0)\|^4 < \infty$$

Note that  $E[U_n] = \mathbf{0}$  and by definition  $U_n(\beta_n^0) = \mathbf{0}$ . Then by Taylor's series expansion we write

$$\begin{aligned} \hat{\beta}_n^0 - \beta^0 &= \mathbf{H}(U_n(\hat{\beta}_n^0)) - \mathbf{H}(U_n(\beta^0)) \\ &= -\frac{\partial \mathbf{H}(U_n)}{\partial U_n} U_n + \frac{1}{2} U_n^T \frac{\partial^2 \mathbf{H}(U_n)}{\partial U_n \partial U_n^T} U_n - \frac{1}{3!} U_n U_n^T \frac{\partial^3 \mathbf{H}(U_n)}{\partial U_n \partial U_n^T \partial U_n} U_n \\ &\quad + \frac{1}{4!} U_n U_n^T \frac{\partial^4 \mathbf{H}(U_n)}{\partial U_n \partial U_n^T \partial U_n \partial U_n^T} U_n U_n^T \\ &\quad + \dots \end{aligned} \tag{3.32}$$

As in (2.1.30) consider the  $l^{\text{th}}$  component of  $\mathbf{H}(U_n)$  say  $h_l(U_n)$  and present the corresponding Taylor's expansion of  $h_l(U_n)$  about it's expected null vector as

$$\begin{aligned} h_l(U_n(\hat{\beta}_n^0)) - h_l(U_n(\beta^0)) &= -\frac{\partial h_l(U_n)}{\partial U_n} U_n + \frac{1}{2} U_n^T \frac{\partial^2 h_l(U_n)}{\partial U_n \partial U_n^T} U_n - \frac{1}{3!} U_n U_n^T \frac{\partial^3 h_l(U_n)}{\partial U_n \partial U_n^T \partial U_n} U_n \\ &\quad + \frac{1}{4!} U_n U_n^T \frac{\partial^4 h_l(U_n)}{\partial U_n \partial U_n^T \partial U_n \partial U_n^T} U_n U_n^T \end{aligned}$$

$$+ \dots \quad (3.33)$$

Taking the expectation of both sides of (3.33) the first term becomes zero. The second term can be written as

$$\begin{aligned} E[U_n^T \frac{\partial^2 h_l(U_n)}{\partial U_n^T \partial U_n} U_n] \\ &= E[\text{Tr}(\underline{\mathbf{A}} \underline{\mathbf{B}})] = E[\sum_{i=1}^k \sum_{j=1}^k a_{ij}^{(l)} b_{ji}] \\ &= \frac{1}{n} [\sum_{i=1}^k \sum_{j=1}^k a_{ij}^{(l)} \sigma_{ji}], \end{aligned} \quad (3.34)$$

where

$$\underline{\mathbf{A}} = \frac{\partial^2 h_l(U_n)}{\partial U_n^T \partial U_n} = (a_{ij}^{(l)})_{i,j=1,\dots,k}$$

$$(a_{ji}^{(l)}) = \frac{\partial^2 h_l(U_n)}{\partial U_{n_j} \partial U_{n_i}}.$$

$$\underline{\mathbf{B}} = U_n U_n^T = (b_{ji}),$$

$$b_{ji} = (U_n)_j (U_n)_i \quad \text{and}$$

$$\sigma_{ji} = E[b_{ji}]$$

$$= \left( -\frac{1}{n} \sum_{i=1}^n (\mathbf{b}'(\mathbf{z}_i, \beta^0))^T \varphi(\mathbf{z}_i, \beta^0) \mathbf{b}'(\mathbf{z}_i, \beta^0) \right)_{ji} \quad \text{is the } (j,i)^{\text{th}} \text{ element of the}$$

covariance matrix of  $U_n(\hat{\beta}_n^0)$ .

The third term in (3.33) is obtained from differentiating (3.34) with respect to  $U_n$  and taking the expectation and is expressed as

$$E\left[\sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k c_r a_{rij}^{(l)} b_{ji}\right] = \frac{1}{n^2} \sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k a_{rij}^{(l)} \gamma_{rji}$$

where

$$a_{rij}^{(l)} = \left( \frac{\partial a_{ij}^{(l)}}{\partial U_{nji}} \right)_r = \frac{\partial^3 h_l(U_n)}{\partial U_{nr} \partial U_{nj} \partial U_{ni}}$$

$$\gamma_{rji} = E[c_r b_{jr}] = E[U_{nr} U_{nj} U_{ni}]$$

$$= \left( -\frac{1}{n} \sum_{i=1}^n \mathbf{b}'(\mathbf{z}_i, \beta^0)^T \mathbf{b}'(\mathbf{z}_i, \beta^0) \Gamma(\mathbf{z}_i, \beta^0) \mathbf{b}'(\mathbf{z}_i, \beta^0) \right)_{rji}$$

Differentiating the third term once more with respect to  $U_n$ , the fourth term in (3.33) is expressed as

$$\begin{aligned} E\left[\sum_{s=1}^k \sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k c_s c_r a_{rij}^{(l)} b_{ji}\right] \\ = \frac{1}{n^2} \left[ \sum_{s=1}^k \sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k a_{srij}^{(l)} \lambda_{srji} \right] \end{aligned}$$

where  $\lambda_{srji}$  is the  $(s r j i)^{th}$  element of the fourth central moment of  $U_n(\hat{\beta}_n^0)$

As in (2.1.34) we can then express (3.33) as

$$\begin{aligned}
E[h_l( U_n(\widehat{\beta}_n^0) ) - h_l( U_n(\beta^0) ) ] = & \frac{\sum_{i=1}^k \sum_{j=1}^k a_{ij}^{(l)} \sigma_{ji}}{n} \\
& + \frac{\sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k a_{rij}^{(l)} \gamma_{rji} + \sum_{s=1}^k \sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k a_{srij}^{(l)} \lambda_{srji}}{n^2} \\
& + O(n^{-3}), \quad l = 1, \dots, k
\end{aligned} \tag{3.35}$$

and hence the bias of  $\widehat{\beta}_n^0$  from (3.32) becomes

$$E[\widehat{\beta}_n^0 - \beta^0] = \frac{1}{n} \mathbf{a}(\beta^0) + \frac{1}{n^2} \mathbf{b}(\beta^0) + o\left(\frac{1}{n^2}\right) \quad \text{as } n \rightarrow \infty, \tag{3.36}$$

where

$$\mathbf{a}(\beta^0) = \left( \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k a_{ij}^{(l)} \sigma_{ji} \right) \text{ and}$$

$$\mathbf{b}(\beta^0) = \left( \frac{1}{3!} \sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k a_{rij}^{(l)} \gamma_{rji} + \frac{1}{4!} \sum_{s=1}^k \sum_{r=1}^k \sum_{i=1}^k \sum_{j=1}^k a_{srij}^{(l)} \lambda_{srji} \right)$$

are functions of  $\beta, \phi$  and  $z_i$ 's independent of  $n$ .

In section 3.2 we shall be concerned with the Jackknife estimation aimed at eliminating the leading term in the bias of  $\widehat{\beta}_n^0$  given by  $\frac{1}{n} \mathbf{a}(\beta^0)$ .

### 3.3 Jackknife Estimator of the Regression Model Parameters

We wish to establish the asymptotic bias and distribution of the Jackknife

estimate  $\tilde{\beta}^J$  of  $\beta^0$ . The Jackknife estimator is obtained by deleting in succession one row of the observation vector  $\mathbf{X}$  and computing the MLE  $\hat{\beta}_{n-1}^{0(i)}$  and is defined below.

Let  $\hat{\beta}_n^0$  and  $\hat{\beta}_{n-1}^{0(i)}$  be the MLE's of  $\beta^0$  based on all the  $n$  and  $n-1$  observations, respectively, and define the pseudo values

$$\hat{\beta}_{n,i}^0 = n \hat{\beta}_n^0 - (n-1) \hat{\beta}_{n-1}^{0(i)}, \quad i = 1, \dots, n \quad (3.37)$$

The Jackknife estimate of  $\beta^0$  is obtained by averaging the pseudo-values (3.37) and is defined as

$$\begin{aligned} \hat{\beta}_n^J &= \frac{1}{n} \sum_{i=1}^n \hat{\beta}_{n,i}^0 \\ &= \hat{\beta}_n^0 - (n-1) \frac{1}{n} \sum_{i=1}^n (\hat{\beta}_{n-1}^{0(i)} - \hat{\beta}_n^0). \end{aligned} \quad (3.38)$$

In order to study the Jackknife estimators of  $\beta^0$  we need establish the properties of the delete-one estimators  $\hat{\beta}_{n-1}^{0(i)}$ , obtained from the estimating equations with one observation deleted from the data. Since the explicit form of  $\hat{\beta}_n^0$  may not be available we rely on the estimating equations  $U_n(\hat{\beta}_n^0)$  and it's Taylor's expansion around  $\beta^0$  to derive an approximate expression for  $\hat{\beta}_n^0$ .

Let  $U_{n-1}^{(i)}(\beta^0)$  denote the delete-one estimating equation for the estimate  $\hat{\beta}_{n-1}^{0(i)}$  of  $\beta^0$  defined as

$$U_{n-1}^{(i)}(\beta^0) = U_n(\beta^0) - (\mathbf{b}'(\mathbf{z}_i, \beta^0))^T [\mathbb{T}(\mathbf{X}_i) + \psi(\mathbf{b}(\mathbf{z}_i, \beta^0))], \quad (3.39)$$



whereby  $\hat{\beta}_{n-1}^{0(i)}$  is the solution to the equation  $U_{n-1}^{(i)}(\beta^0) = 0$ . Note that when evaluated at  $\hat{\beta}_n^0$ , (3.39) satisfies

$$U_{n-1}^{(i)}(\hat{\beta}_n^0) = -(\mathbf{b}'(\mathbf{z}_i, \beta^0))^T \left\{ \mathbb{T}(X_i) + \psi(\mathbf{b}(\mathbf{z}_i, \beta^0)) \right\} \Big|_{\hat{\beta}_n^0}. \quad (3.40)$$

From (3.40) one obtains the following

$$U_{n-1}^{(i)}(\hat{\beta}_n^0) - U_{n-1}^{(i)}(\hat{\beta}_{n-1}^{0(i)}) = -(\mathbf{b}'(\mathbf{z}_i, \beta^0))^T \left\{ \mathbb{T}(X_i) + \psi(\mathbf{b}(\mathbf{z}_i, \beta^0)) \right\} \Big|_{\hat{\beta}_n^0}. \quad (3.41)$$

Unlike in the non regression estimation problem considered in Chapter 2, (3.41) the terms  $\mathbb{T}(X_i)$  appear as a product with functions of the estimates of  $\beta^0$ . So that we need a different approach, from the one used in Chapter 2, to study the behavior of the difference in the estimates  $(\hat{\beta}_{n-1}^{0(i)} - \hat{\beta}_n^0)$ , the delete-one estimates and the estimates from the entire sample.

To examine the behavior of  $\sum_{i=1}^n \{\hat{\beta}_{n-1}^{0(i)} - \hat{\beta}_n^0\}$  in (3.38) we need expressions of  $\hat{\beta}_{n-1}^{0(i)}$  and  $\hat{\beta}_n^0$ . Since the explicit expressions are often not available in various estimation problems, we derive an approximate representation of  $\hat{\beta}_n^0$  by considering the Taylor's expansion of  $U_n(\hat{\beta}_n^0)$ . It will be helpful at this time to introduce the following notations and results from matrix algebra which will be used in studying the approximate expression for  $\hat{\beta}_n^0$ .

Let

$$\mathbf{v}_i(\beta^0) = -(\mathbf{b}'(\mathbf{z}_i, \beta^0))^T \varphi(\mathbf{b}(\mathbf{z}_i, \beta^0)) \mathbf{b}'(\mathbf{z}_i, \beta^0), \quad i = 1, \dots, n,$$

$$\mathbf{V}_n(\beta^0) = \sum_{i=1}^n \mathbf{v}_i(\beta^0),$$

$$\mathbf{u}_i(\boldsymbol{\beta}^0) = (\mathbf{b}'(\mathbf{z}_i, \boldsymbol{\beta}^0))^T [\mathbf{T}(X_i) + \boldsymbol{\psi}(\mathbf{b}(\mathbf{z}_i, \boldsymbol{\beta}^0))],$$

$$\mathbf{U}_n(\boldsymbol{\beta}^0) = \sum_{i=1}^n \mathbf{u}_i(\boldsymbol{\beta}^0).$$

So that for the delete-one expressions we have

$$\begin{aligned} \mathbf{V}_{n-1}^{(i)}(\boldsymbol{\beta}^0) &= \mathbf{V}_n(\boldsymbol{\beta}^0) - \mathbf{v}_i(\boldsymbol{\beta}^0) \text{ and} \\ \mathbf{U}_{n-1}^{(i)}(\boldsymbol{\beta}^0) &= \mathbf{U}_n(\boldsymbol{\beta}^0) - \mathbf{u}_i(\boldsymbol{\beta}^0), \end{aligned} \quad i = 1, \dots, n.$$

The matrix  $\mathbf{V}_n(\boldsymbol{\beta}^0)$  is symmetric and positive definite and note that  $\frac{1}{n}\mathbf{V}_n(\boldsymbol{\beta}^0)$  converges to  $\bar{\mathbf{I}}_n(\boldsymbol{\beta}^0)$  which is the average Fisher information. Furthermore assuming each of the  $\mathbf{v}_i(\boldsymbol{\beta}^0)$  is small relative to  $\mathbf{V}_n(\boldsymbol{\beta}^0)$  we can obtain the  $(\mathbf{V}_{n-1}^{(i)}(\boldsymbol{\beta}^0))^{-1}$ , the inverse of  $\mathbf{V}_{n-1}^{(i)}(\boldsymbol{\beta}^0)$  as

$$\begin{aligned} (\mathbf{V}_{n-1}^{(i)}(\boldsymbol{\beta}^0))^{-1} &= (\mathbf{V}_n(\boldsymbol{\beta}^0) - \mathbf{v}_i(\boldsymbol{\beta}^0))^{-1} \\ &= \{ \mathbf{I} + (\mathbf{V}_n(\boldsymbol{\beta}^0))^{-1} \mathbf{v}_i(\boldsymbol{\beta}^0) + (\mathbf{V}_n(\boldsymbol{\beta}^0))^{-1} \mathbf{v}_i(\boldsymbol{\beta}^0) \mathbf{v}_i(\boldsymbol{\beta}^0)^T (\mathbf{V}_n(\boldsymbol{\beta}^0))^{-1} + \dots \} (\mathbf{V}_n(\boldsymbol{\beta}^0))^{-1} \\ &= (\mathbf{V}_n(\boldsymbol{\beta}^0))^{-1} + (\mathbf{V}_n(\boldsymbol{\beta}^0))^{-1} \mathbf{v}_i(\boldsymbol{\beta}^0) (\mathbf{V}_n(\boldsymbol{\beta}^0))^{-1} + O(n^{-3}). \end{aligned}$$

Also the following result can be proved by matrix algebra,

$$(\mathbf{V}_n(\boldsymbol{\beta}^0))^{-1} - (\mathbf{V}_{n-1}^{(i)}(\boldsymbol{\beta}^0))^{-1} = (\mathbf{V}_n(\boldsymbol{\beta}^0))^{-1} - \{ \mathbf{I} - (\mathbf{V}_n(\boldsymbol{\beta}^0))^{-1} \mathbf{v}_i(\boldsymbol{\beta}^0) \}^{-1} (\mathbf{V}_n(\boldsymbol{\beta}^0))^{-1}$$

$$= (\mathbf{V}_n(\beta^0))^{-1} \mathbf{v}_i(\beta^0) (\mathbf{V}_n(\beta^0))^{-1} + O(n^{-3}).$$

Now the Taylor's expansion of  $U_n(\hat{\beta}_n^0)$  around  $\beta^0$  given by

$$\begin{aligned} U_n(\hat{\beta}_n^0) &= U_n(\beta^0) + \frac{\partial U_n(\beta^0)}{\partial \beta^0} (\hat{\beta}_n^0 - \beta^0) + \frac{1}{2} (\hat{\beta}_n^0 - \beta^0)^T \frac{\partial^2 U_n(\beta^0)}{\partial \beta^0 \partial \beta^0} (\hat{\beta}_n^0 - \beta^0) \\ &+ \dots \end{aligned} \quad (3.43)$$

and multiplying both sides of (3.43) by  $(\mathbf{V}_n(\beta^0))^{-1}$  and noting that  $U_n(\hat{\beta}_n^0) = 0$ , for sufficiently large  $n$ ,  $\partial U_n(\beta^0) / \partial \beta^0 \rightarrow \mathbf{V}_n(\beta^0)$ , we obtain

$$\hat{\beta}_n^0 - \beta^0 = (\mathbf{V}_n(\beta^0))^{-1} \{U_n(\beta^0) + W_n(\beta^0)\},$$

where

$W_n(\beta^0) = \frac{1}{2} (\hat{\beta}_n^0 - \beta^0)^T \frac{\partial^2 U_n(\beta^0)}{\partial \beta^0 \partial \beta^0} (\hat{\beta}_n^0 - \beta^0) = O(1)$  and similarly for the delete-one estimating equations,

$$\hat{\beta}_{n-1}^{0(i)} - \beta^0 = (\mathbf{V}_{n-1}^{(i)}(\beta^0))^{-1} \{U_{n-1}^{(i)}(\beta^0) + \mathbf{W}_{n-1}^{(i)}(\beta^0)\}. \quad (3.44)$$

By substituting the expressions for  $(\mathbf{V}_{n-1}^{(i)}(\beta^0))^{-1}$  and  $U_{n-1}^{(i)}(\beta^0)$  given above we obtain, for large  $n$ ,

$$\begin{aligned} \hat{\beta}_n^0 - \hat{\beta}_{n-1}^{0(i)} &= (\mathbf{V}_n(\beta^0))^{-1} \mathbf{u}_i(\beta^0) - (\mathbf{V}_n(\beta^0))^{-1} \mathbf{v}_i(\beta^0) (\mathbf{V}_n(\beta^0))^{-1} U_n(\beta^0) \\ &+ (\mathbf{V}_n(\beta^0))^{-1} \mathbf{v}_i(\beta^0) (\mathbf{V}_n(\beta^0))^{-1} \mathbf{u}_i(\beta^0) - (\mathbf{V}_n(\beta^0))^{-1} \mathbf{v}_i(\beta^0) (\mathbf{V}_n(\beta^0))^{-1} O(1) \end{aligned}$$

$$+ O(n^{-3}).$$

If we write  $\mathbf{W}_{n-1}^{(i)}(\beta^0) = \mathbf{W}_n(\beta^0) - (\mathbf{W}_n(\beta^0) - \mathbf{W}_{n-1}^{(i)}(\beta^0))$ , then

$$\sum_{i=1}^n \hat{\beta}_n^0 - \hat{\beta}_{n-1}^{0(i)} = \sum_{i=1}^n (\mathbf{V}_n(\beta^0))^{-1} \mathbf{v}_i(\beta^0) (\mathbf{V}_n(\beta^0))^{-1} (\mathbf{u}_i(\beta^0) - \mathbf{W}_n(\beta^0)) + O(n^{-2}). \quad (3.45)$$

Since  $\mathbf{u}_i(\beta^0) = O_p(1)$  and  $(\mathbf{V}_n(\beta^0))^{-1} = O(n^{-1})$ , the first term on the RHS of (3.45) has a negligible contribution to the left hand side of (3.45) and the second term is  $(\mathbf{V}_n(\beta^0))^{-1} O_p(1) = O_p(n^{-1}) = o_p(n^{-1/2})$ .

Hence  $\sum_{i=1}^n \hat{\beta}_n^0 - \hat{\beta}_{n-1}^{0(i)} = o_p(n^{-1/2})$ , which implies that

$$\sqrt{n} \sum_{i=1}^n \{ \hat{\beta}_{n-1}^{0(i)} - \hat{\beta}_n^0 \} \xrightarrow{p} 0, \quad \text{as } n \rightarrow \infty. \quad (3.46)$$

So that from (3.38) we have

$$\sqrt{n}(\hat{\beta}_n^{0J} - \hat{\beta}_n^0) = o_p(1).$$

Since  $\hat{\beta}_n^0$  is asymptotically normal by (3.29) we conclude by Slutsky's theorem that for large  $n$

$$\sqrt{n}(\hat{\beta}_n^{0J} - \beta^0) \xrightarrow{\mathcal{D}} \mathbf{N}(0, \bar{\mathbf{I}}_n^{-1}(\beta^0)). \quad (3.47)$$

For inference purposes then, the Jackknifed estimate could be used in place of the MLE.

To demonstrate the reduction in bias of the MLE  $\hat{\beta}_n^0$  achieved by jackknifing, we

express the bias of the Jackknife estimator as

$$\begin{aligned}
 E[\tilde{\beta}^J - \beta^0] &= E[\hat{\beta}_n^0 - \beta^0] - \frac{(n-1)}{n} \sum_{i=1}^n E[(\hat{\beta}_{n-1}^{0(i)} - \hat{\beta}_n^0)] \\
 &= E[\hat{\beta}_n^0 - \beta^0] - \frac{(n-1)}{n} \sum_{i=1}^n E[(\hat{\beta}_{n-1}^0 - \hat{\beta}_n^0)] \\
 &= nE[\hat{\beta}_n^0 - \beta^0] - (n-1)E[(\hat{\beta}_{n-1}^0 - \hat{\beta}_n^0)]
 \end{aligned}$$

and using the results obtained earlier (3.36) on the bias of  $\hat{\beta}_n^0$  we have

$$E[\tilde{\beta}^J - \beta^0] = \frac{1}{n(n-1)} \mathbf{b}(\beta^0) + o\left(\frac{1}{n^2}\right) \quad \text{as } n \rightarrow \infty$$

where  $\mathbf{b}(\beta^0)$  is as defined in (3.36).

## Chapter 4

# Length-Biased Distribution for Exponential Families.

### 4.1 Introduction.

Let  $X$  be a continuous r.v. following a length-biased distribution with p.d.f.,  $g(X; \theta)$  which is a weighted p.d.f. of an exponential family of distributions given by

$$g(X; \theta) = \frac{X f(X; \theta)}{\gamma(\theta)} = c^*(X) \exp\{ \underline{T}(X)^T \theta + \Psi^*(\theta) \}, \quad X > 0 \quad (4.1)$$

where

$\theta = (\theta_1, \dots, \theta_k)$  is a  $k$ -parameter vector

$$\Psi^*(\theta) = \Psi(\theta) - \log \gamma(\theta)$$

$$\gamma(\theta) = E_f[X]$$

$\underline{T}(X) = (T_1(X_i), \dots, T_k(X_i))^T$  is a  $k$ -vector of sufficient statistics for  $\theta$

$\Psi^*(\cdot)$  and  $\gamma(\cdot)$  are known functions of their arguments.

The density  $g(X; \theta)$  is said to be a length-biased version of the pdf  $f(X; \theta)$ . We study the bias and asymptotic distribution of the MLE and jackknifed MLE of  $\theta$  obtained from the sample of size  $n$  drawn from the distribution  $g(X; \theta)$ . Note that the densities  $f(X; \theta)$  and  $g(X; \theta)$  are both of exponential family of distributions,

hence with appropriate assumptions on  $\Psi^*(\theta)$ , the results we have for the density  $f(X;\theta)$  apply as well to the density  $g(X;\theta)$ . Also the parameter  $\theta$  is the same for the two distributions, however as will be illustrated in section 4.2 the estimates of  $\theta$  under the different likelihoods will differ.

## 4.2 Maximum Likelihood Estimators of the Parameters

Following similar development discussed in Chapter 2, we have, under the density  $g(X;\theta)$ ,

$$\begin{aligned} E_g[\mathcal{T}(X)] &= -\psi^*(\theta) \\ &= -\psi(\theta) + \zeta'(\theta) \end{aligned}$$

$$\text{where } \zeta'(\theta) = \frac{\partial \log \gamma(\theta)}{\partial \theta} = \frac{\gamma'(\theta)}{\gamma(\theta)}$$

This clearly indicates the sample mean  $\bar{\mathcal{T}}_n(X)$  from the length-biased data will be a biased estimator of the original population mean from the density  $f(X;\theta)$ . Furthermore, in a univariate setup with  $\mathcal{T}(X_i) = X_i$  the bias is the ratio of the variance to the mean of  $X_i$ , which is a positive quantity.

Next we consider the estimation of the parameter vector  $\theta$  under the length-biased distribution (4.1). For the maximum likelihood estimation of  $\theta$ , we assume the same set of assumptions "B" as in Chapter 2 with  $\Psi(\theta)$  replaced by  $\Psi^*(\theta)$ . The estimating equations for  $\theta$  are then given by

$$U_n(\boldsymbol{\theta}) = \sum_{i=1}^n \frac{\partial \log g(X_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{i=1}^n \mathfrak{T}(X_i) + n \boldsymbol{\psi}^*(\boldsymbol{\theta}) = \mathbf{0}. \quad (4.2)$$

So that the MLE  $\hat{\boldsymbol{\theta}}_n$  of  $\boldsymbol{\theta}$  is the solution for  $\boldsymbol{\theta}$  in (4.2). The gradient of the estimating equation is obtained by differentiating (4.2) with respect to  $\boldsymbol{\theta}^T$  and is expressed as

$$-E_{\boldsymbol{\theta}} \frac{\partial U_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^T} = -n\boldsymbol{\varphi}^*(\boldsymbol{\theta}), \quad \text{finite and p.d.} \quad (4.3)$$

Note that the MLE  $\hat{\boldsymbol{\theta}}_n$  obtained from solving (4.2) for  $\boldsymbol{\theta}$ , will be different from the MLE obtained from the original unweighted distribution and the the two estimates can not be compared unless their explicit expressions are available, even then the expressions may not be simple enough for any general comparison to be made. This will be true for comparing the bias of the two estimates, however since the expected Fisher's information from the two probability models is a function of only the parameter  $\boldsymbol{\theta}$  conditions under which one model is more informative than the other will be investigated.

### 4.3 Asymptotic Distribution and Bias of the MLE

First we establish the asymptotic normal distribution of  $\hat{\boldsymbol{\theta}}_n$  from (4.2). As noted above the difference in the length-biased version of the exponential family of distributions (4.1) from the original distribution is in the reparametrization of the function  $\Psi(\cdot)$ , under the regularity set of assumptions "B" we conclude using the same development, that



$$\sqrt{n}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) \xrightarrow{\mathfrak{D}} \mathcal{N}(\mathbf{0}, \mathbf{I}_g^{-1}(\boldsymbol{\theta})), \quad (4.4)$$

where  $\mathbf{I}_g(\boldsymbol{\theta}) = -\boldsymbol{\varphi}^*(\boldsymbol{\theta})$ , the subscript  $g$  denotes the length-biased distribution, is the Fisher information from the length-biased distribution characterized by  $g(\mathbf{X};\boldsymbol{\theta})$ . In section 4.4 we compare the Fisher information from the density  $f(\mathbf{X};\boldsymbol{\theta})$  with the Fisher information from  $g(\mathbf{X};\boldsymbol{\theta})$ .

The structure of the bias of  $\widehat{\boldsymbol{\theta}}_n$ , obtained from (4.2) is similar to the form of bias studied in Chapter 2 except for the quantities involved. So that as in (2.1.35) the leading term in the bias of  $\widehat{\boldsymbol{\theta}}_n$  under the distribution  $g(\mathbf{X};\boldsymbol{\theta})$  is given as

$$\frac{1}{n}\mathbf{a}^*(\boldsymbol{\theta}) = \left( \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k a_{ij}^{(l)*} \sigma_{ji}^* \right),$$

which is a similar expression for the leading term in the bias of  $\widehat{\boldsymbol{\theta}}_n$  under  $f(\mathbf{X};\boldsymbol{\theta})$  with  $a_{ij}^{(l)}$  and  $\sigma_{ji}$  replaced by  $a_{ij}^{(l)*}$  and  $\sigma_{ji}^*$  respectively.

In particular, as in proposition 1, let  $\boldsymbol{\eta}^* = -\boldsymbol{\psi}^*(\boldsymbol{\theta})$  and a function  $\mathbf{H}: \mathbb{R}^k \rightarrow \mathbb{R}^k$  given by  $\mathbf{H}^*(\boldsymbol{\eta}^*) = \boldsymbol{\theta}$ , with  $\mathbf{H}(\boldsymbol{\eta}^*) = (h_1^*(\boldsymbol{\eta}^*), \dots, h_l^*(\boldsymbol{\eta}^*), \dots, h_k^*(\boldsymbol{\eta}^*))^T$ . Define the sample mean

$$\bar{\mathbf{T}}_n = (\bar{\mathbf{T}}_{n1}, \dots, \bar{\mathbf{T}}_{nk}), \text{ where}$$

$$\bar{\mathbf{T}}_{nj} = \frac{1}{n} \sum_{i=1}^n \mathbf{T}_{nj}(\mathbf{X}_i) \quad j = 1, \dots, k.$$

The  $\sigma_{jj}^*$  are the elements of the covariance matrix of  $\bar{\mathbf{T}}_n$  based on the length-biased distribution and are given by

$$\sigma_{jj'}^* = \mathbb{E}_g[\bar{\mathbf{T}}_{nj} - \eta_j^*][\bar{\mathbf{T}}_{nj'} - \eta_{j'}^*] \quad j, j' = 1, \dots, k$$

$$\begin{aligned}
&= E_g[\bar{T}_{nj} - \eta_j - \zeta'_j(\boldsymbol{\theta})][\bar{T}_{nj'} - \eta_{j'} - \zeta'_{j'}(\boldsymbol{\theta})] \\
&= \sigma_{jj'} - (\zeta'(\boldsymbol{\theta})\zeta'(\boldsymbol{\theta})^T)_{jj'}
\end{aligned}$$

where  $\zeta'_j(\boldsymbol{\theta})$  and  $(\zeta'(\boldsymbol{\theta})\zeta'(\boldsymbol{\theta})^T)_{jj'}$  are the  $j^{\text{th}}$  and  $jj'^{\text{th}}$  elements of  $\zeta'(\boldsymbol{\theta})$  and  $(\zeta'(\boldsymbol{\theta})\zeta'(\boldsymbol{\theta})^T)$  respectively and  $\sigma_{jj'}$  are the  $jj'$  elements of the covariance matrix of  $\bar{\mathbf{T}}_n$  under the original probability distribution  $f$ .

$$a_{ij}^{(l)*} = \frac{\partial^2 \mathbf{h}_l^*(\boldsymbol{\eta}^*)}{\partial \eta_j^* \partial \eta_i^*}$$

Recall that by definition  $\mathbf{H}^*(\boldsymbol{\eta}^*) = \boldsymbol{\theta}$ , and  $\mathbf{H}(\boldsymbol{\eta}) = \boldsymbol{\theta}$ , where  $\boldsymbol{\eta}^* = \boldsymbol{\eta} + \zeta'(\boldsymbol{\theta})$  so that if we let  $\zeta_{jj'} = (\zeta'(\boldsymbol{\theta})\zeta'(\boldsymbol{\theta})^T)_{jj'}$  we observe that

$$\sum_{i=1}^k \sum_{j=1}^k a_{ij}^{(l)*} \sigma_{ji}^* = \sum_{i=1}^k \sum_{j=1}^k a_{ij}^{(l)*} (\sigma_{ji} - \zeta_{ji}).$$

This shows the leading term in the bias of the length-biased distribution involves both the bias from the original distribution and the bias induced by the length-biased design. Next we compare the information about  $\boldsymbol{\theta}$  provided by the estimates  $\hat{\boldsymbol{\theta}}_n$  from the samples drawn from  $f(\mathbf{X};\boldsymbol{\theta})$  and  $g(\mathbf{X};\boldsymbol{\theta})$ .

The goal here is to investigate theoretically which of the two designs provide greater information about the unknown parameters. Fisher information is used to compare the efficiency for the parameters in two design models under which samples are obtained by comparing their respective optimality. Optimal design criteria is defined by a function  $\Phi_p$  defined on the matrix  $\mathbf{C}$ , where  $\Phi_p$  satisfies the

following conditions.

- i)  $\Phi_p$  is nondecreasing and invariant with respect to row and column permutations.
- ii) If  $C_1 \geq C_2$ , i.e.  $C_1 - C_2$  is non negative definite, then

$$\Phi(C_1) \leq \Phi(C_2)$$

Keifer (1975) defines some of the widely used optimality criteria as;

$$\Phi_p(C) = \left( \frac{1}{q-1} \sum \lambda^{-p} \right)^{\frac{1}{p}},$$

where  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_{q-1}$  are the eigenvalues of matrix C.

Silvey, 1980; Shah, 1989 discussed optimality criteria measures applicable in the comparison of any two information matrices  $C_1$  and  $C_2$ . If  $C_1 - C_2$  is nonnegative definite matrix then  $\Phi_p(C_1) \geq \Phi_p(C_2)$  where  $\Phi_p(C)$  is a measure of the amount of information in C. The three  $\psi$  measures proposed are for values of  $p = 0, 1$ , and  $\infty$ , in  $\Phi_p$  corresponding to D-optimality, A-optimality and E-optimality respectively. If C is an inverse of the Fisher information matrix for  $\theta$ , D-optimality minimizes  $\det(C)$ , provides a measure of the generalized variance of the components of  $\hat{\theta}$ . A-optimality minimizes the  $\text{trace}(C)$  and represents the average variance of the elements of  $\hat{\theta}$ . E-optimality minimizes the largest variance of the components of  $\hat{\theta}$  and is a measure of the largest variance of the components of  $\hat{\theta}$ .

From (4.2) we obtain the average Fisher information about  $\theta$  as

$$\begin{aligned} I_g(\theta) &= E\left[ \frac{1}{n} \frac{\partial U_n(\theta)}{\partial \theta^T} \right] = -\varphi^*(\theta) \\ &= -\varphi(\theta) + \zeta''(\theta) \\ &= I_f(\theta) + \zeta''(\theta), \quad \text{finite and p.d.} \end{aligned} \tag{4.5}$$

From (4.5) we have the following implications

$$|\mathbf{I}_g(\boldsymbol{\theta})| = |\mathbf{I}_f(\boldsymbol{\theta}) + \zeta''(\boldsymbol{\theta})| \leq |\mathbf{I}_f(\boldsymbol{\theta})| + |\zeta''(\boldsymbol{\theta})| \quad (4.6)$$

where  $|A|$  denotes determinant of  $A$ .

So that

$|\mathbf{I}_g(\boldsymbol{\theta})| \leq |\mathbf{I}_f(\boldsymbol{\theta})|$  if  $|\zeta''(\boldsymbol{\theta})| \leq 0$ , that is, if  $\zeta''(\boldsymbol{\theta})$  is negative semi-definite. The equality of the two information matrices will be realized if  $\zeta''(\boldsymbol{\theta}) \equiv \mathbf{0}$ .

On the other hand if  $\zeta''(\boldsymbol{\theta})$  is positive definite then  $|\zeta''(\boldsymbol{\theta})| > 0$ , and D-optimality criteria is indeterminate.

The trace of (4.5) is

$$\text{tr}(\mathbf{I}_g(\boldsymbol{\theta})) \quad \begin{cases} < \text{tr}(\mathbf{I}_f(\boldsymbol{\theta})) & \text{if } \text{tr}(\zeta''(\boldsymbol{\theta})) < 0, \\ > \text{tr}(\mathbf{I}_f(\boldsymbol{\theta})) & \text{if } \text{tr}(\zeta''(\boldsymbol{\theta})) > 0 \end{cases}$$

which implies that by E-optimality criteria, the length-biased design would be more informative than the original design if  $\text{tr}(\zeta''(\boldsymbol{\theta})) > 0$  and otherwise sampling from the original distribution would be more informative about  $\boldsymbol{\theta}$ .

The following remarks are also noteworthy in making the comparisons.

### Remarks

- 1) If  $\zeta(\boldsymbol{\theta})$  is a polynomial in  $\boldsymbol{\theta}$ , of degree  $r < 2$ , then  $\zeta''(\boldsymbol{\theta}) \equiv \mathbf{0}$ , in which case the information about  $\boldsymbol{\theta}$  obtained from the length-biased design and the original design is equivalent.
- 2) For  $k = 1$  and  $\zeta(\boldsymbol{\theta})$  is a convex function of  $\boldsymbol{\theta}$ , then  $\zeta''(\boldsymbol{\theta}) > 0$ , so that the length-biased design is more informative about  $\boldsymbol{\theta}$  than the original design.

**Example 4.1** Let  $X_i$ ,  $i = 1, \dots, n$ , be i.i.d.r.v.'s following a lognormal distribution  $(\mu, \sigma^2)$  as given in example 1 of Chapter 2. It is easy to show that  $E_f(X_i) = \exp(\mu + \frac{\sigma^2}{2})$ . Recall that  $\theta_1 = \frac{\mu}{\sigma^2}$ ,  $\theta_2 = -\frac{1}{2\sigma^2}$  and so that  $\zeta(\theta) = \log E_f[X] = -\frac{1}{2\theta_2}(\theta_1 + \frac{1}{2\theta_2})$  and

$$\Psi(\theta) = \frac{1}{2} \left( \frac{\theta_1^2}{2\theta_2} - \log \frac{-1}{2\theta_2} \right).$$

This implies that

$$\varphi(\theta) = \frac{1}{2} \begin{bmatrix} \frac{1}{\theta_2} & -\frac{\theta_1}{\theta_2^2} \\ -\frac{\theta_1}{\theta_2^2} & \frac{\theta_1^2}{\theta_2^3} + \frac{1}{\theta_2^2} \end{bmatrix} \quad \text{and} \quad \zeta''(\theta) = \frac{1}{2} \begin{bmatrix} 0 & \frac{1}{\theta_2^2} \\ \frac{1}{\theta_2^2} & -\frac{2\theta_1}{\theta_2^3} - \frac{1}{\theta_2^3} \end{bmatrix}. \quad (4.7)$$

Since  $|\zeta''(\theta)| = -\frac{1}{4\theta_2^4} < 0$ , by D-optimality we conclude that length-biased sampling design is less informative than sampling from the original original distribution for the lognormal probability model, in the sense that the generalized variance from the length-biased design is greater than the generalized variance from sampling from the original distribution.

The trace of  $\zeta''(\theta)$  is  $\text{tr}(\zeta''(\theta)) = -\left(\frac{\theta_1}{\theta_2^3} + \frac{1}{2\theta_2^3}\right) > 0$ , since  $\theta_1 > 0$  and  $\theta_2 < 0$ , so that by A-optimality the length-biased design would be desirable since it provides the minimum average variance of the components of  $\hat{\theta}$ .

#### 4.4. Estimating Regression Parameters from the Length-Biased Distribution

Let  $X_i$  ( $i = 1, \dots, n$ ) be a r.v. following a pdf  $g(X_i; \beta, \phi)$ , which is a length-biased version of the distribution of  $f(X_i; \beta, \phi)$  in (3.2), given by

$$g(X_i; \beta, \phi) = \frac{X_i f(X_i; \beta, \phi)}{\gamma(\mathbf{b}(z_i, \beta, \phi))} = c^*(X_i) \exp\left\{ \mathbf{b}(z_i, \beta, \phi)^T \underline{T}(X_i) + \Psi^*(\mathbf{b}(z_i, \beta, \phi)) \right\}, \quad X_i > 0$$

where  $\Psi^*(\mathbf{b}(z_i, \beta, \phi)) = \Psi(\mathbf{b}(z_i, \beta, \phi)) - \log \gamma(\mathbf{b}(z_i, \beta, \phi))$  and

$\gamma(\mathbf{b}(z_i, \beta, \phi)) = E_f[X_i] < \infty$ , where  $E_f$  denotes expectation with respect to the density  $f(\cdot)$ .

Recall that  $\underline{T}(X_i) = (T_1(X_i), \dots, T_q(X_i))^T$ , so that we may find a one-to-one transformation  $h$ , such that

$$h(T_j(X_i)) = X_i \text{ and } E_f|h(T_j(X_i))| < \infty, \text{ for some } 1 \leq j \leq k$$

This is certainly the case with the lognormal distribution with  $T_1(X_i) = \log(X_i)$ , and  $\mathbf{b}_1(z_i, \beta, \sigma^2) = z_i^T \beta / \sigma^2$ .

As before we derive the moments of  $\underline{T}(X_i)$  from length-biased density as

$$E_g[\underline{T}(X_i)] = -\psi^*(\mathbf{b}(z_i, \beta^0)).$$

The second central moment of  $\mathbb{T}(X_i)$  is given by

$$\begin{aligned} E_g[\mathbb{T}(X_i) - \psi^*(\mathbf{b}(z_i, \beta^0))][\mathbb{T}(X_i) - \psi^*(\mathbf{b}(z_i, \beta^0))]^T \\ = -\varphi^*(\mathbf{b}(z_i, \beta^0)). \end{aligned}$$

Expressing the moments of  $\mathbb{T}(X_i)$  in terms of the moments of the original density we have

$$E_g[\mathbb{T}(X_i)] = -\left\{ \psi(\mathbf{b}(z_i, \beta^0)) - \frac{\gamma'(\mathbf{b}(z_i, \beta^0))}{\gamma(\mathbf{b}(z_i, \beta^0))} \right\}.$$

As an illustration if  $q = 1$ ,  $T_1(X_i) = X_i$ . (The gamma distribution  $G(\lambda, \nu)$  with  $\nu$  known and  $\gamma = z^T \beta$  would fit this case).

In this case,

$$E_f[X_i] = -\psi(\mathbf{b}(z_i, \beta^0)) = \gamma(\mathbf{b}(z_i, \beta^0))$$

and

$$\begin{aligned} E_g[X_i] &= -\left\{ \psi(\mathbf{b}(z_i, \beta^0)) + \frac{\varphi(\mathbf{b}(z_i, \beta^0))}{\psi(\mathbf{b}(z_i, \beta^0))} \right\} \\ &= E_f[X_i] + \frac{\text{Var}_f(X_i)}{E_f[X_i]}. \end{aligned}$$

This implies the sample mean from the length-biased distribution would be a biased estimate of the population mean of the true distribution. To obtain the MLE of  $\beta$ , and  $\phi$ , of the length-biased data from  $g(X_i; \beta, \phi)$  we maximize the log-likelihood function in the same fashion as before. So the MLE  $\hat{\beta}_n$ , and  $\hat{\phi}_n$  of

$\beta$ , and  $\phi$ , are the solutions to the equations

$$U_n(\beta^0) = \sum_{i=1}^n (\mathbf{b}'(\mathbf{z}_i, \beta^0))^T [\mathbb{T}(X_i) + \psi^*(\mathbf{b}(\mathbf{z}_i, \beta^0))] = \mathbf{0}. \quad (4.8)$$

Note that the estimating equation is not a sum of identically distributed r.v.'s, however the terms are independent, hence we can still use some large sample results to establish the asymptotic normality of the MLE. The estimating equations (4.8) are similar to the estimating equations (3.12) and by similar arguments leading to (3.25) we conclude  $U_n(\beta^0)$  in (4.7) has an asymptotic normal distribution

$$\frac{1}{\sqrt{n}} U_n(\beta^0) \longrightarrow^{\mathcal{D}} \mathcal{N}(\mathbf{0}, \mathbf{I}_g(\beta^0)), \quad (4.9)$$

where  $\mathbf{I}_g(\beta^0)$  is the limiting average Fisher information from the length-biased likelihood given below. Also in line with (3.29) we have the following result,

$$\sqrt{n}(\hat{\beta}_n^0 - \beta^0) \longrightarrow^{\mathcal{D}} \mathcal{N}(\mathbf{0}, \mathbf{I}_g^{-1}(\beta^0)) \quad (4.10)$$

The gradient of the estimating equation  $U_n(\beta^0)$  is given by

$$\begin{aligned} \frac{\partial U_n(\beta^0)}{\partial \beta^{0T}} &= \sum_{i=1}^n (\mathbf{b}''(\mathbf{z}_i, \beta^0))^T [\mathbb{T}(X_i) + \psi^*(\mathbf{b}(\mathbf{z}_i, \beta^0))] \\ &\quad + \sum_{i=1}^n (\mathbf{b}'(\mathbf{z}_i, \beta^0))^T \varphi^*(\mathbf{b}(\mathbf{z}_i, \beta^0))]^T \mathbf{b}'(\mathbf{z}_i, \beta^0). \end{aligned}$$

For every  $i = 1, \dots, n$  we define



$$\begin{aligned} I_i(\beta^0) &= E_g \left[ \frac{\partial \log g(X_i; \beta^0)}{\partial \beta^0} \right] \left[ \frac{\partial \log g(X_i; \beta^0)}{\partial \beta^0} \right]^T \\ &= -(\mathbf{b}'(z_i, \beta^0))^T \varphi^*(\mathbf{b}(z_i, \beta^0)) \mathbf{b}'(z_i, \beta^0), \end{aligned}$$

so that, as  $n \rightarrow \infty$  the average Fisher information is

$$\begin{aligned} \bar{I}_g(\beta^0) &= -\frac{1}{n} E \left[ \frac{\partial U_n(\beta^0)}{\partial \beta^{0T}} \right] \\ &= -\frac{1}{n} \sum_{i=1}^n (\mathbf{b}'(z_i, \beta^0))^T \varphi^*(\mathbf{b}(z_i, \beta^0)) \mathbf{b}'(z_i, \beta^0). \end{aligned}$$

The average Fisher information matrix for  $\beta^0$ , from the length-biased distribution can be expressed in terms of the Fisher information for  $\beta^0$  from the original distribution by substituting in the expression for  $\varphi^*(\mathbf{b}(z_i, \beta^0))$ , to obtain

$$\bar{I}_g(\beta^0) = \bar{I}_f(\beta^0) + \zeta''(\beta^0),$$

$$\text{where } \zeta''(\beta^0) = \frac{1}{n} \sum_{i=1}^n (\mathbf{b}'(z_i, \beta^0))^T \left\{ \frac{\gamma''(\mathbf{b}(z_i, \beta^0))}{\gamma(\mathbf{b}(z_i, \beta^0))} - \frac{\gamma'(\mathbf{b}(z_i, \beta^0)) \gamma'(\mathbf{b}(z_i, \beta^0))^T}{(\gamma(\mathbf{b}(z_i, \beta^0)))^2} \right\} \mathbf{b}'(z_i, \beta^0).$$

so that the Fisher information from the length-biased distribution will be less than or greater or equal to the Fisher information depending on the nature of  $\log \gamma(\mathbf{b}(z_i, \beta^0))$ .

Remarks:

(1). If  $\log \gamma(\mathbf{b}(z_i, \beta^0))$  is a polynomial of degree  $r$ , ( $r \leq 1$ ), in  $\mathbf{b}(z_i, \beta^0)$ , then  $\zeta''(\beta^0)$

vanishes. In this case we get the same amount of information from the two probability models.

(2). If  $q = 1$ , and  $\log\gamma(b(z_i, \beta^0))$  is a convex function of  $b(z_i, \beta^0)$  then the length-biased design is more informative about  $\beta^0$  than the original design.

# Chapter 5

## Example and Simulations

### 5.1 Introduction

The previous chapters studied the asymptotic bias reduction and optimality of the designs. Since the asymptotics tells us how the estimates behave in large samples, it is important to study also how well the proposed estimators perform in small-to-moderate sample sizes. In this chapter we study applications to a real data problem which is of length-biased design, and also study the performance of our estimators in moderate sample sizes.

The data set for analysis in this exercise is from the 1989/1990 Sudan Demographic and Health Surveys (DHS). Our objective is to estimate the mean birth interval between the first and second births among married women, 14-49 years old, who have been married for at least 10 years. Without loss of generality we assume the women constitutes a cohort married at the same time ( $t = 0$ ) and with complete birth intervals, that is, both the first and second births dates are available. The data then consist of  $N = 3211$  complete birth intervals. In addition to the birth intervals, the covariates age, duration of marriage use for the women were also extracted for the regression analysis.

In order to obtain a length-biased sample of all the complete birth intervals, we select the birth intervals straddling the 5-year period (survey time,  $\tau = 5$ ) for the sample. Such a sampling scheme is known to induce a length-biased sample

since longer birth intervals are more likely to overlap any arbitrarily chosen time, and hence be sampled, than shorter birth intervals. The length-biased sample selected consists of  $n = 681$  birth intervals. From the 681 birth intervals we obtained  $M = 1000$  sub-samples of size  $n = 50, 100$  and  $200$  to study and compare the performance of the jackknife estimates and the MLE. The estimates are the averages of the 1000 estimates from each of the  $n = 50, 100$  and  $200$  sample sizes.

We assume the population of  $N = 3211$  birth intervals follow a lognormal distribution with the baseline survival function given by

$$S_W(w) = 1 - \Phi((w - \mu)/\sigma), \quad (5.1)$$

where  $w$  is the natural logarithm of the birth intervals,  $\mu$  and  $\sigma$  are the location and shape parameters, respectively.  $\Phi$  is the cumulative distribution function of the normal distribution.

The lognormal regression model with the covariates age (months) and duration of marriage (months) is fit to the data to estimate biases and the estimation of the effects of the covariate terms. The model assumed is

$$W = \beta_0 + \beta_1 Z_1 + \beta_2 Z_2 + \sigma \epsilon, \quad (5.2)$$

where

$W$  is the natural logarithm of the birth interval,

$\beta_0, \beta_1$  and  $\beta_2$  are the regression coefficients,

$\sigma$  is scale parameter,

$Z_1$  and  $Z_2$  denote Age and Duration of marriage for the women, respectively, and

$\epsilon$  is the error term.

The baseline survival function is of the form (5.1) with  $\mu = \beta_0 + \beta_1 Z_1 + \beta_2 Z_2$

For the model (5.1) without covariates, we denote the parameter vector by  $\theta = (\mu, \sigma)$  and its MLE's and jackknife estimates are denoted by  $\hat{\theta}_l$  and  $\hat{\theta}_l^J$ , respectively,  $l = 1, \dots, M$ . The sample mean and variance of  $\theta$  computed from all the  $M$  replicate samples are

$$\hat{\theta} = \frac{1}{M} \sum_{l=1}^M \hat{\theta}_l$$

and

$$\hat{V} = \frac{1}{M-1} \sum_{l=1}^M (\hat{\theta}_l - \hat{\theta})(\hat{\theta}_l - \hat{\theta})^T.$$

The standard error of the components of  $\hat{\theta}$  is estimated as the square root of the corresponding diagonal entry of  $\hat{V}$ .

The overall performance of our estimators is shown by the root mean squared error (RMSE). The RMSE measures the combined effect of bias and sampling variation, and for the parameter  $\mu$  the expected value of the MSE is equal to:

$$E(\hat{\mu} - \mu)^2 = \sigma^2(\hat{\mu}) + (E(\hat{\mu}) - \mu)^2. \quad (5.3)$$

The RMSE for the estimate  $\hat{\mu}$  of  $\mu$  is estimated by

$$\text{RMSE} = \sqrt{\frac{1}{1000} \sum_{l=1}^{1000} (\hat{\mu}_l - \mu)^2}. \quad (5.4)$$

The RMSE of  $\hat{\sigma}^2$  is defined similarly, with  $\hat{\mu}$  in (5.4) replaced by  $\hat{\sigma}^2$ .

In addition to the parameter estimates, the distribution of the M estimates from the replicate samples is provided to support our methodological claim that the jackknifed estimates and MLE are asymptotically normal. The distributions are presented graphically for the parameters ( $\mu$  and  $\sigma$ ) by way of both the smoothed density function obtained from the empirical distribution of the estimates.

## 5.2 Results and Discussion

The parameters for the model (5.1) for the  $N = 681$  birth intervals are  $\mu = 3.5984$  and  $\sigma = 0.5550$ . The results of our estimates from fitting model (5.3) to the Sudan birth interval data are given in Table 5.1.

Table 5.1  
Mean, Bias and RMSE of the MLE and Jackknife Estimates  
for  $M = 1000$  samples of size  $n$  from a Lognormal Distribution Model.  
1990 Sudan DHS birth interval data.

$n$	MLE						Jackknife Estimates					
	$\hat{\mu}$	Bias	RMSE	$\hat{\sigma}$	Bias	RMSE	$\hat{\mu}$	Bias	RMSE	$\sigma$	Bias	RMSE
50	3.6005	0.0021	0.4179	0.5458	-0.0092	0.0720	3.6005	0.0021	0.4179	0.5546	-0.0004	0.0782
100	3.5986	0.0002	0.4122	0.5518	-0.0032	0.0632	3.5986	0.0002	0.4122	0.5562	0.0012	0.0669
200	3.5989	0.0005	0.4103	0.5540	-0.0010	0.0570	3.5989	0.0005	0.4103	0.5562	0.0012	0.0590
681	3.5984	-		0.5550	-							

## Discussion

Both the MLE and jackknifed estimates of the location parameter  $\mu$  were identical suggesting the leading term in the bias of MLE for  $\mu$  is zero, hence the classical jackknifing did not improve the MLE. The jackknife estimate for the scale parameter, for each of the sample sizes 50, 100, and 200 is less biased for the length-biased estimate of the  $n = 681$  compared to the corresponding MLE estimate. The MLE of the scale parameter performed better than the jackknife estimate in terms of a smaller RMSE, however, the difference in the RMSE of the two estimates is small. This lack of gain in the jackknifed estimator over the MLE for the location parameter, in this example, of the lognormal model (5.1), is explained by the lack of contribution of the leading term in the bias of the MLE. Recall that the jackknifed estimator reduces the bias of order  $n^{-1}$  from the bias of the MLE and for large samples this reduction is minimal. The leading term in the bias of the MLE  $\hat{\theta}_n$  is a function  $\mathbf{a}(\theta)$  whose components are linear combinations of the third partial derivatives of the loglikelihood. In the case of the model (5.1) the coordinate of  $\mathbf{a}(\theta)$  for  $\mu$  vanishes asymptotically, resulting in no gain in the bias reduction of the Jackknife estimator. To see this, recall from Chapter 4 that the first element of  $\Psi(\theta)$  is  $\theta_1^2/\theta_2$ , so that the  $(1,1)^{th}$  element of  $\Gamma(\theta)$  is zero.

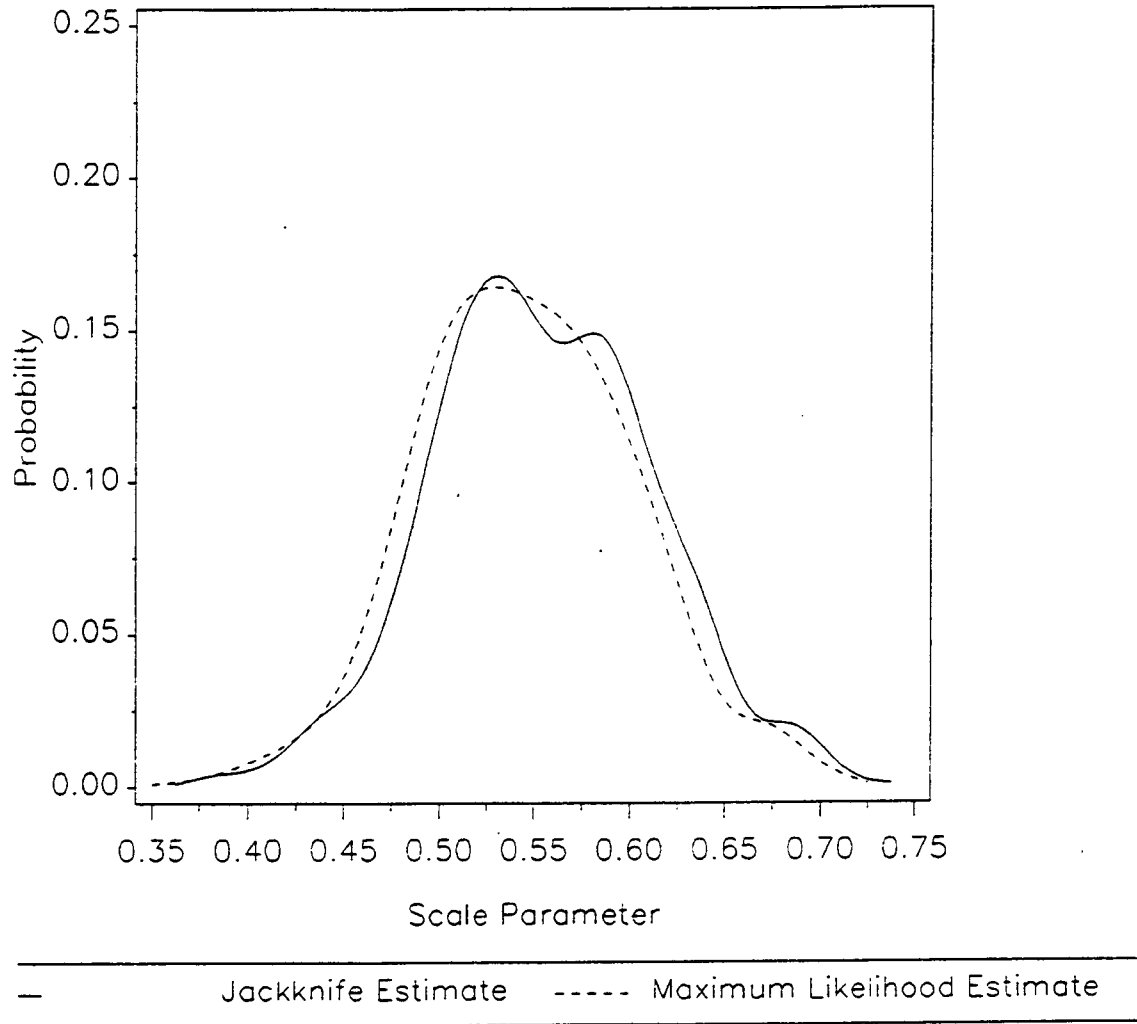
Figure 5.1 shows the distribution of the 1000 estimates of the scale parameter for the samples of size  $n = 100$ . The smoothed density plots shows a good approximation to the normal distribution for the sample size of  $n = 100$ . The distribution of the shape parameter  $\mu$  for the same samples (not provided here) also exhibit a good approximation to the normal distribution. The same result, of no difference in the MLE and the jackknifed estimate for the location parameter is

also shown by the coincidence of the density plots for the two estimates of  $\mu$ . The distribution of the estimates from the sample sizes  $n(=50,200)$  also showed a good approximation to the normal distribution.



# Figure 5.1

Maximum Likelihood and Jackknifed Estimates  
of the Scale Parameter of the Lognormal Distribution  
for the SUDAN First Order Birth Interval Data  
SUDAN DHS (1990)



## Regression Model

The lognormal regression model with covariates, duration of marriage and age of the woman, is fit to the length-biased sample from the Sudan birth interval data. The length-biased sample consists of  $n = 681$  birth intervals. The population regression coefficients are  $\beta_0 = 3.8228$ ,  $\beta_1 = -0.0006$ ,  $\beta_2 = -0.0002$  and  $\sigma = 0.5514$ , where  $\beta_0$ ,  $\beta_1$  and  $\beta_2$  are the intercept duration and age coefficients respectively. From the  $n=681$  we obtained  $M = 1000$  samples of size  $n=50$ ,  $100$  and  $200$  to study and compare the performance of our the jackknife estimator and the MLE. The estimates are the averages from the  $1000$  samples for each sample size. The parameter vector for the regression model is denoted by  $\boldsymbol{\beta} = (\beta_0, \beta_1, \beta_2, \sigma)$  and it's corresponding MLE and jackknife estimates are obtained in the same fashion as discussed above for estimating the parameter  $\boldsymbol{\theta}$ .

The results are shown in Table 5.2. The largest bias reduction is observed for the scale parameter for all sample sizes, where the MLE have a negative bias and the jackknifed estimates are biased upwards. The RMSE's from the MLE and the jackknifed MLE for the intercept term and the regression coefficients are comparable, whereas the MLE performs better than the jackknifed MLE for the scale parameter in term of smaller RMSE.

Table 5.2

Mean, Bias and RMSE of the MLE and Jackknife Estimates  
of estimates of mean birth interval from the Lognormal Regression Model,  
1990 Sudan DHS birth interval data.  
M = 1000 samples of size  $n = 50, 100, 200$

$n$		MLE			Jackknife Estimates		
		Mean	Bias	RMSE	Mean	Bias	RMSE
50	$\beta_0$	3.8083	-0.0145	0.7141	3.8037	-0.0191	0.0798
	$\beta_1$	-0.0007	-0.0001	0.0022	-0.0007	-0.0001	0.0022
	$\beta_2$	-0.0001	0.0001	0.0020	-0.0001	0.0001	0.0020
	$\sigma$	0.5321	-0.0193	0.0658	0.5543	0.0029	0.0787
100	$\beta_0$	3.8211	-0.0017	0.6147	3.8176	-0.0052	0.6117
	$\beta_1$	-0.0006	0.0000	0.0016	-0.0006	0.0000	0.0016
	$\beta_2$	-0.0002	0.0000	0.0013	-0.0001	0.0001	0.0013
	$\sigma$	0.5443	-0.0071	0.0579	0.5543	0.0029	0.0659
200	$\beta_0$	3.8233	0.0005	0.5651	3.8254	0.0026	0.5630
	$\beta_1$	-0.0006	0.0000	0.0009	-0.0006	0.0000	0.0009
	$\beta_2$	-0.0002	0.0000	0.0087	-0.0002	0.0001	0.0013
	$\sigma$	0.5484	-0.0030	0.0529	0.5534	0.0020	0.0574

$N = 681$  with true parameters  $\beta_0 = 3.8228$ ,  $\beta_1 = -0.0006$ ,  $\beta_2 = -0.0002$  and  $\sigma = 0.5514$

The density estimates of the estimates for the  $M = 1000$  samples of size  $n = 50$  are shown in Figures 5.2, 5.3, 5.4 and 5.5 for the intercept term, Age and Duration regression coefficients and the scale parameters, respectively. The graphs illustrates similar approximately normal distributions for both the MLE and the Jackknife estimates. The graphs of the estimates from the other sample sizes (not shown here) showed similar patterns.

Figure 5.2

Maximum Likelihood and jackknifed Estimates of the Intercept Parameter of the Lognormal Regression Model for the SUDAAN First Order Birth Interval Data SUDAAN DHS (1990)

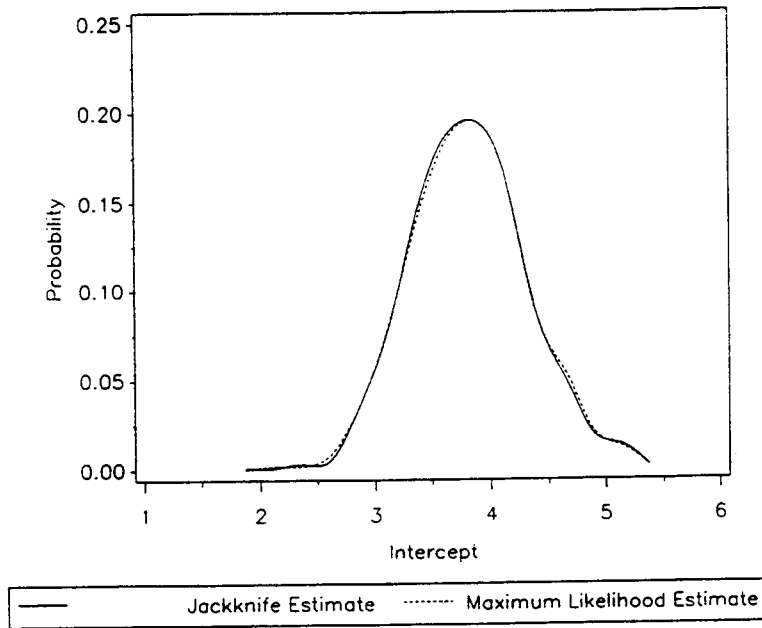


Figure 5.3

Maximum Likelihood and jackknifed Estimates of the Marital Duration Parameter of the Lognormal Regression Model for the SUDAAN First Order Birth Interval Data SUDAAN DHS (1990)

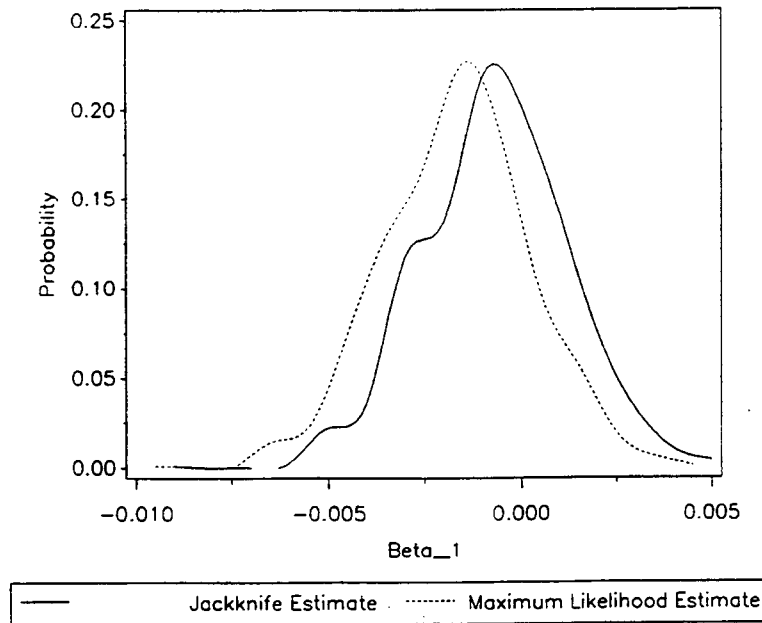


Figure 5.4

Maximum Likelihood and jackknifed Estimates of the Age Parameter of the Lognormal Regression Model for the SUDAAN First Order Birth Interval Data SUDAAN DHS (1990)

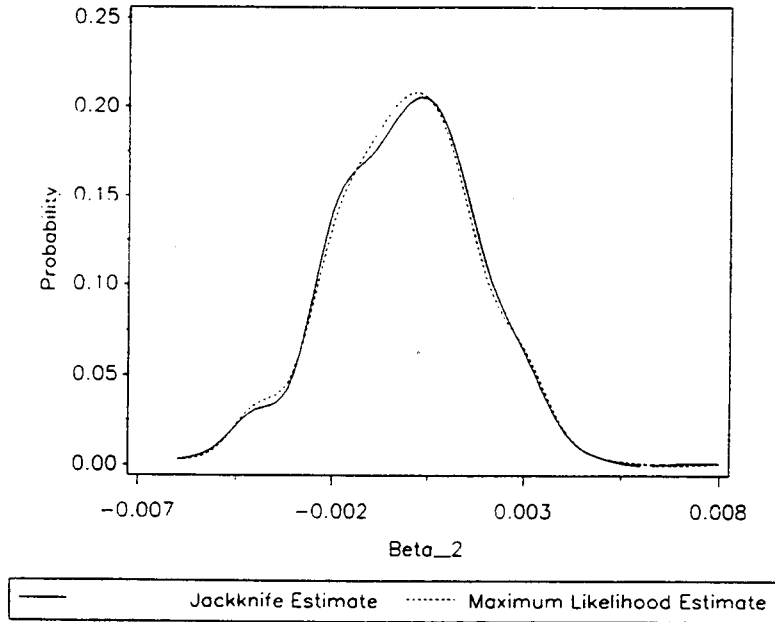
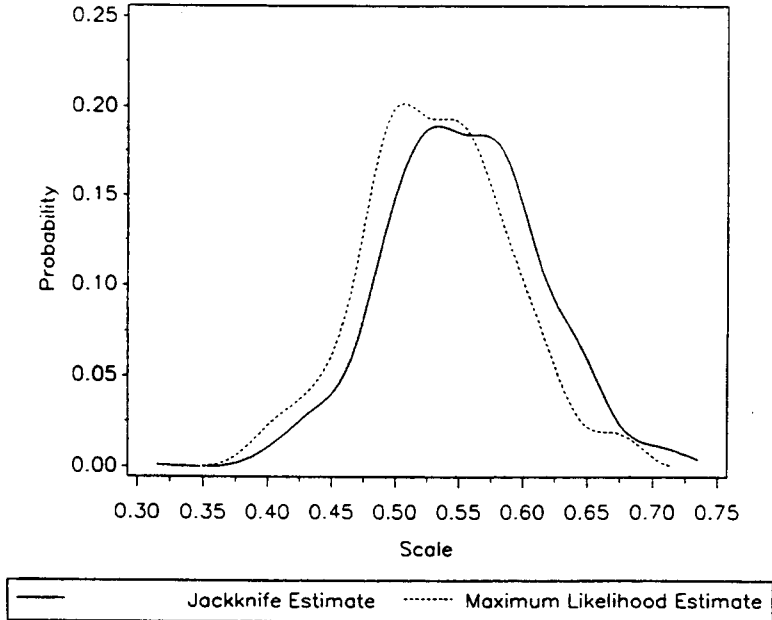


Figure 5.5

Maximum Likelihood and jackknifed Estimates of the Scale Parameter of the Lognormal Regression Model for the SUDAAN First Order Birth Interval Data SUDAAN DHS (1990)



### 5.3 Simulation Study

The performance of our estimates in moderate samples is studied through simulation studies. The main focus of our study is on the performance of the jackknife estimates compared to the MLE. Two distributions for the response variate are chosen for this study, and these are the lognormal and the gamma. For each distribution we consider estimating the parameters of the model without the covariate terms. The regression analysis is studied through simulation study where the response variable follows a lognormal distribution. In each case our simulations involve  $m=1000$  data sets, each containing samples of size  $n = 50, 100, \text{ and } 200$ .

Our first study concerns the estimation of the parameters  $\mu$  and  $\sigma$  and obtains their root mean squared errors. In particular, we obtain samples from the following distributions:

(1) the lognormal, with location parameter  $\mu = 0.5$  and scale parameter  $\sigma = 1.5$ ; and

(2) the gamma with location parameter  $\mu = 1.5$  and  $\sigma^2 = 1.5$ .

The baseline survival functions for the lognormal and gamma distribution models with parameters shape parameter  $\mu$  and scale parameter  $\sigma^2$  are given by:

(1) Lognormal:

$$S_W(w) = 1 - \Phi((w - \mu)/\sigma); \text{ and} \tag{5.5}$$

(2) Gamma:

$$S_X(x) = \int_x^\infty \frac{\sigma}{\Gamma(\mu)} (\sigma t)^{\mu-1} \exp\{-\sigma t\} dt, \quad (5.6)$$

where  $w$  is the natural logarithm of the response variate and  $\Phi$  is the cumulative distribution function of the normal distribution.

It is easy to show that the length-biased version of the original distribution from the lognormal and the gamma distributions preserves their original distributional forms. The length-biased distribution with the original form being the lognormal with parameters  $\mu$  and  $\sigma^2$  is itself lognormal with the same parameter  $\sigma^2$  and a shape parameter  $\mu^* = \mu + \sigma^2$ . For the gamma distribution the shape parameter for the length-biased density becomes  $\mu^* = \mu + 1$ , where  $\mu$  is the shape parameter of the original distribution gamma and the scale parameter  $\sigma$  remains unchanged.

The lognormal random variate with mean  $e^{\{\mu + \sigma^2/2\}}$  and variance  $e^{\{2\mu + \sigma^2\}}(e^{\sigma^2} - 1)$  is generated using the SAS system's function NORMAL. The NORMAL function returns a random variate generated from a normal distribution with mean zero and variance 1. To generate the random variate following the gamma distribution with the shape parameter  $\mu$  and scale parameter  $\sigma$ , we use the SAS RANGAM function. The RANGAM function returns a random variate distributed according to the gamma function with the specified parameter values.

The first simulation samples were drawn from the lognormal model with  $\mu = 2.0$  and  $\sigma^2 = 1.5$ . The simulations from the gamma were obtained with  $\mu = 1.5$  and  $\sigma^2 = 1.5$ . For each distribution  $M = 1000$  samples of sizes  $n = 50, 100$  and  $200$  were obtained. The MLE and the jackknifed estimates are shown in Tables 5.3 and 5.4 for the Lognormal and Gamma distributions, respectively.



Table 5.3  
 Mean, Bias and RMSE of the MLE and Jackknife Estimates  
 for  $M = 1000$  samples of size  $n$  from a Lognormal Distribution.  
 Simulations with  $\mu = 2.0$  and  $\sigma^2 = 1.5$

$n$	MLE						Jackknife Estimates					
	$\hat{\mu}$	Bias	RMSE	$\hat{\sigma}$	Bias	RMSE	$\hat{\mu}$	Bias	RMSE	$\sigma$	Bias	RMSE
50	1.9997	-0.0003	0.1795	1.2013	-0.0023	0.3220	1.9997	-0.0003	0.1795	1.2198	-0.0049	0.3057
100	1.9993	-0.0007	0.1244	1.2122	-0.0013	0.3000	1.9993	-0.0007	0.1244	1.2213	0.0034	0.2914
200	1.9983	-0.0017	0.0849	1.2214	-0.0033	0.2864	1.9983	-0.0017	0.0849	1.2260	0.0013	0.2809

The simulation results from the lognormal distribution confirm the findings in the real data example that there is no significant bias reduction by jackknifing, for the MLE of the shape parameter. However jackknifing had a positive result in reducing the bias of the MLE for the scale parameter.

Table 5.4  
 Mean, Bias and RMSE of the MLE and Jackknife Estimates  
 for  $M = 1000$  samples of size  $n$  from a Gamma Distribution.  
 Simulations with  $\mu = 1.5$  and  $\sigma^2 = 1.5$

$n$	MLE						Jackknife Estimates					
	$\hat{\mu}$	Bias	RMSE	$\hat{\sigma}$	Bias	RMSE	$\hat{\mu}$	Bias	RMSE	$\sigma$	Bias	RMSE
50	1.1601	-0.3399	0.3640	0.8703	-0.3544	0.3437	1.1668	-0.3332	0.3577	0.8853	-0.3394	0.3538
100	1.1654	-0.3328	0.3475	0.8805	-0.3459	0.3368	1.1688	0.3295	0.3443	0.8880	-0.3383	0.3441
200	1.1654	-0.3346	0.3407	0.8816	-0.3431	0.3380	1.1670	-0.3330	0.3391	0.8853	-0.3394	0.3428

The simulations from the gamma probability distribution indicate a sizable reduction in the bias of the MLE by jackknifing. The largest bias reduction by

jackknifing is observed for the scale parameter. The RMSE of the jackknifed estimate is slightly larger than the RMSE of the MLE of the scale parameter. So the bias reduction by jackknifing the MLE of the scale parameter is achieved at the expense of having a larger RMSE.

Next we consider fitting the regression model to each of the model lognormal model (5.5). The covariates  $Z_1$  and  $Z_2$  are generated from bivariate normal with parameters  $\mu_1 = 0$ ,  $\mu_2 = 0$ ,  $\sigma_1 = 1$ ,  $\sigma_2 = 1$ ,  $\rho = 0.35$ . The simulation study from the lognormal distributions should provide us with a typical applications for the proposed estimates of the length-biased data from exponential family of distributions.

For the linear regression model we assume

$$Y = \beta_0 + \beta_1 Z_1 + \beta_2 Z_2 + \sigma \epsilon, \quad (5.8)$$

where  $\epsilon$  have p.d.f.  $f(\epsilon)$ ,  $\beta_j$  ( $j = 0, 1, 2$ ) are unknown regression coefficients to be estimated together with the scale parameter  $\sigma$ .

The sample sizes used in the simulations are  $n=50$ , 100, and 200. In each case we generate 1,000 samples of size  $n$  observations and compute  $\hat{\beta}_j$ . The root mean squared error, is the criterion used to assess the performance of the estimates and is given by

$$\text{RMSE} = \sqrt{\frac{1}{1000} \sum_{i=1}^{1000} (\hat{\beta}_{ji} - \beta_j)^2} \quad (j = 0, 1, 2)$$

for the regression model parameter estimates. The results are shown in Table 5.5. The MLE and the jackknifed estimate for the intercept term are similar.

Jackknifing improved the bias of the MLE of the scale parameter. With respect to the RMSE the jackknifed estimates performed better than the MLE of the scale parameter, whereas for the intercept and regression coefficients the RMSE for the MLE and the jackknifed estimate are comparable.

Table 5.5

Mean, Bias and RMSE of the MLE and Jackknife Estimates of the Parameters of the Lognormal Regression Model,  $M = 1000$  simulations of samples of size  $n = 50, 100, 200$

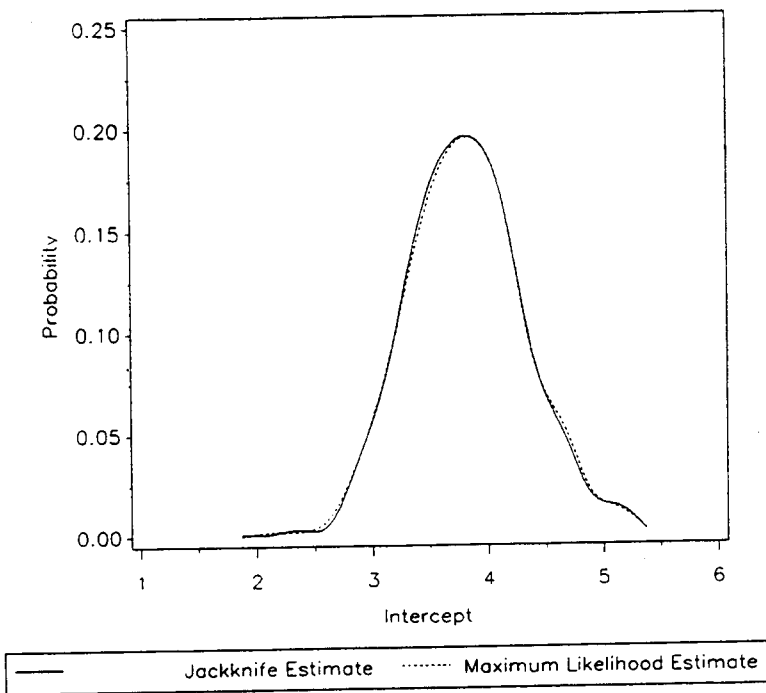
$n$		MLE			Jackknife Estimates		
		Mean	Bias	RMSE	Mean	Bias	RMSE
50	$\beta_0$	2.2496	0.7405	0.7698	2.2496	0.7405	0.7698
	$\beta_1$	1.4921	-0.0079	0.1853	1.4915	-0.0085	0.1870
	$\beta_2$	1.8041	0.0041	0.2027	1.8041	0.0041	0.2031
	$\sigma$	1.1685	-0.0563	0.1394	1.2125	-0.0122	0.1329
100	$\beta_0$	2.2506	0.7506	0.7605	2.2506	0.7506	0.7605
	$\beta_1$	1.4936	-0.0064	0.1289	1.4933	-0.0067	0.1291
	$\beta_2$	1.8048	0.0048	0.1398	1.8047	0.0047	0.1398
	$\sigma$	1.2005	-0.0247	0.0913	1.2222	-0.0025	0.0896
200	$\beta_0$	2.2502	0.7502	0.7546	2.2502	0.7502	0.7546
	$\beta_1$	1.4919	-0.0081	0.0897	1.4918	-0.0082	0.0896
	$\beta_2$	1.8053	0.0053	0.0955	1.8052	0.0052	0.0954
	$\sigma$	1.2092	-0.0155	0.0637	1.2201	-0.0046	0.0625

True parameters  $\beta_0 = 1.5$ ,  $\beta_1 = 1.5$ ,  $\beta_2 = 1.8$  and  $\sigma^2 = 1.5$

The distribution of the parameter estimates from the  $M = 1000$  samples of size  $n = 50$  are shown in Figures 5.6, 5.7, 5.8 and 5.9 for the intercept term,  $X_1$  and  $X_2$  regression coefficients and the scale parameters, respectively. The graphs suggest that the distribution of both the MLE and the jackknifed MLE are approximately normal. Figure 5.9 further illustrate the shift in the scale distribution resulting from a sizable bias reduction by jackknifing the MLE of the scale parameter.

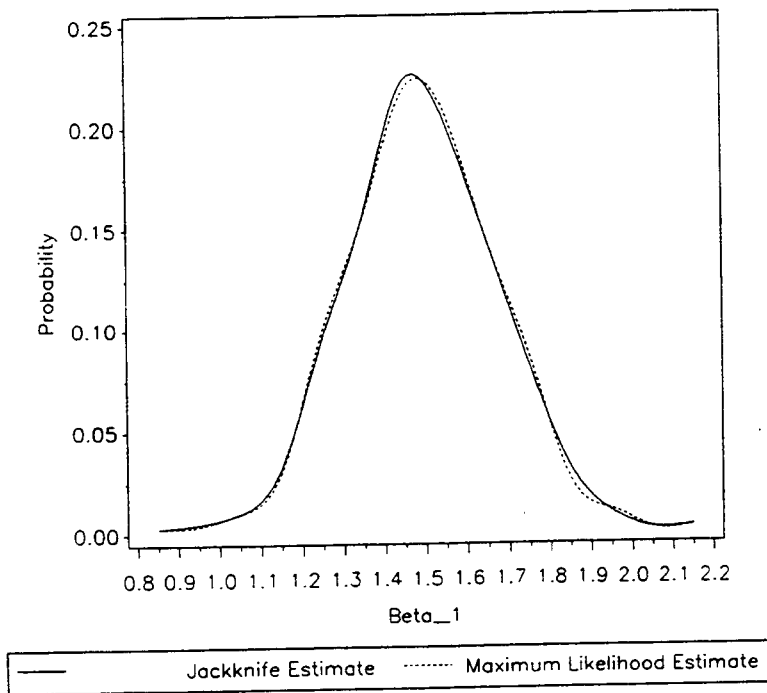
### Figure 5.6

Maximum Likelihood and Jackknifed Estimates of the Intercept Parameter of the Lognormal Regression Model  
The mean estimates from M=1000 samples of size n=50



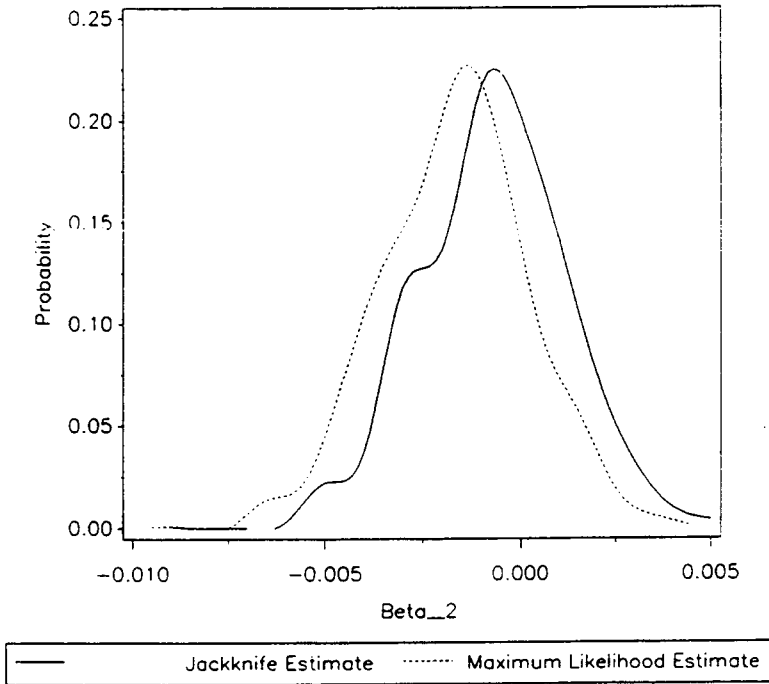
### Figure 5.7

Maximum Likelihood and Jackknifed Estimates of the Beta\_1 Parameter of the Lognormal Regression Model  
The mean estimates from M=1000 samples of size n=50



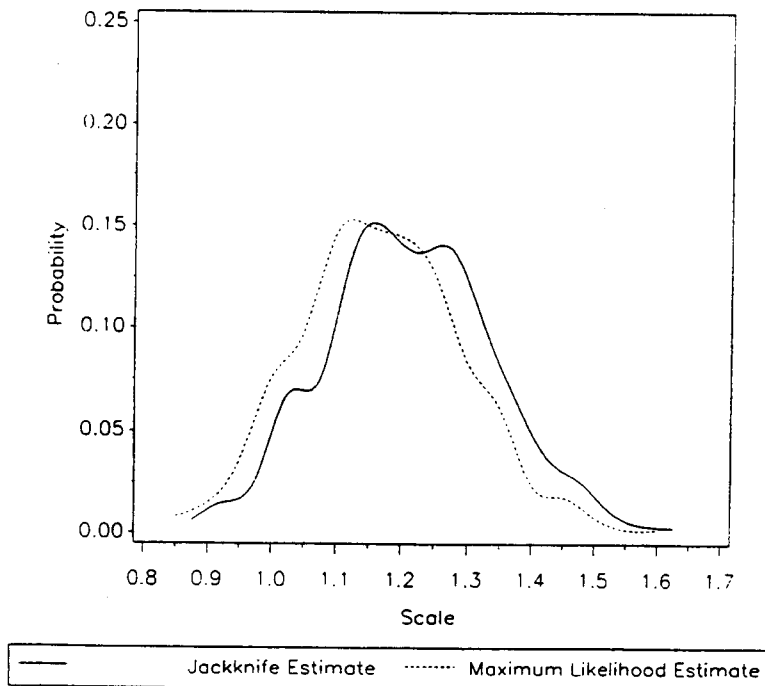
### Figure 5.8

Maximum Likelihood and Jackknifed Estimates of the Beta\_2 Parameter of the Lognormal Regression Model  
The mean estimates from M=1000 samples of size n=50



### Figure 5.9

Maximum Likelihood and Jackknifed Estimates of the Scale Parameter of the Lognormal Regression Model  
The mean estimates from M=1000 samples of size n=50



# Chapter 6

## Discussion and Suggestions for Future Research

The main goal of this research work has been to develop the form of bias of the maximum likelihood estimators for the parameters of the length-biased distribution from an exponential family, and jackknife the MLE to reduce their bias. The classical jackknife estimation technique eliminates the bias of order  $n^{-1}$ , in the MLE. Both the MLE and the Jackknifed estimators are shown analytically to have the same asymptotic normal distribution.

The Fisher information matrix from the original distribution and its length-biased version are compared by D-optimality and A-optimality criteria. The Fisher information from the length-biased density  $g$  of an exponential family is the sum of the Fisher information from the original density  $f$  and a matrix of the second partial derivative, with respect to the parameter vector of the  $\log E_f[Y]$ . It is shown that there are situations when the D-optimality criterion may be inconclusive in determining which design provides more information about the parameters of the model.

Numerical studies are carried out to evaluate the theoretical results given in this research work. The results indicate the MLE are generally biased. However,

in the cases studied, especially for the lognormal distribution, the reduction in the bias of the MLE by jackknifing is minimal. In particular MLE and jackknifed estimates for the shape parameter are identical, indicating no bias of order  $n^{-1}$  in the MLE. The jackknife technique performed as well as the MLE in terms of the RMSE. The results also suggest both the MLE and the Jackknifed estimators have the same asymptotic normal distribution, as claimed in our analytical work.

There are issues for further research to be resolved. This work assumes prior knowledge of the distributional form of the response variable. There are situations when this prior information may not be available. In this case some other estimation techniques, such as semi-parametric or other nonparametric approaches, would be the appropriate, and their jackknifed estimators can be analyzed. Some other family of distributions, other than the exponential family for the original distribution would provide more insight into the study of the optimality of the length-biased design compared to sampling from the original design.

Another area of interest for future research would be to consider the length-biased distribution in the case of multivariate response. In this case marginal length-biased distributions may be developed.



## BIBLIOGRAPHY

- Bayarri, M.J. and De Groot, M.H. (1989). Comparison of Experiments with Weighted Distributions. *Statistical Data Analysis and Inference* (ed Yadolah Dodge), North Holland Publishing Co.,185-196.
- Bayarri, M.J. and De Groot, M.H. (1987). Information in selection models. In *Probability and Bayesian Statistics* (ed. R.Viertl ), 39-51.
- Bhattacharyya, B.B. Franklin, L.A. and Richardson, G.D. (1988) A comparison of nonparametric unweighted and length-biased density of fibres. *Comm. Statist. A 17* , 3629-44.
- Blackwell, D. (1951). Comparison of experiments. *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability*,93-102. Berkeley, California: University of California Press.
- Brillinger, D.R. (1964). The Asymptotic behavior of Tukey's general method of setting approximate confidence limits ( the jackknife ) when applied to maximum likelihood estimates. *Review of the International Statistical Institute. 32*, 202-206.
- Cox, D.R. (1969). Some sampling problems in technology. In *New Developments in Survey Sampling*, (ed. N.L. Johnson and H. Smith), pp. 506-27. New York: Wiley.
- Cramer, H. (1946). *Mathematical Methods of Statistics*. Princeton, NJ: Princeton University Press.
- Godambe, V.P. and Rajarhi, M.B. (1989). Optimal Estimation for Weighted Distributions: Semi-Parametric Model. *Statistical Data Analysis and Inference* (ed Yadolah Dodge), North Holland Publishing Co.,199-208.
- Gray, H.L., Watkins, T.A. and Adams, J.E. (1972). On the jackknife statistic, its extensions and its relation to  $e_n$ -transformations. *Ann. Math. Statist. 43*, 1-30.
- Gupta, R.C. and Keating, J.P. (1986). Relations for Reliability Measures Under Length-Biased Sampling. *Scand. J. Statist. 13*, 49-56 .
- Gupta, R.C. and Kirmani S.N.U.A. (1990). The role of weighted distributions in stochastic modeling. *Comm. Statist. Theory math., 19*, 3147-3162.
- Huber, P.J. (1964). Robust estimation of a location parameter. *Ann. math. and statist., 35*, 73-101.
- Jones, M.C. (1991). Kernel density estimation for length biased data. *Biometrics, 78*, 3, pp. 511-519.
- Kendal, M.G. and Stuart, A. (1958). *The Advanced Theory of Statistics*, Vol. 2 (1958). Charles Griffin, London.

- Krieger, A.M. and Pfeiffermann, D. (1992). Maximum Likelihood Estimation from Complex Sample Surveys. *Survey Methodology*, 18, 225-239.
- Lindley, D.V. (1956). On measure of Information provided by an experiment. *Ann. Math. Statist.*, 27, 968-1005.
- LeCam, L. (1956). On the asymptotic theory of estimation and testing hypotheses. In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*, vol, 1, 129-156. Berkeley: University of California Press.
- McCullagh, P. and Nelder, J.A. (1989). *Generalized Linear Models*, 2nd edition. London; Chapman and Hall.
- Miller, R.G. (1974). The Jackknife - a review. *Biometrika*, 61, 1, 1-15.
- Miller, R.G. (1974). An unbalanced Jackknife. *Ann. Statist.*, 2, 880-891.
- Patil, G.P. and Rao, C.R. (1978). Weighted distributions and size-biased sampling with applications to wildlife populations and human families. *Biometrics* 34, 179-89 .
- Patil, G.P. and Rao, C.R. (1977). Weighted distributions: A Survey of Their Applications. In *Applications of Statistics*, (ed. P.R. Krishnaiah), North Holland Publishing Co., 383-405.
- Rao, C.R. (1985). Weighted Distributions Arising out of Methods of Ascertainment: What Population Does a Sample represent? *A celebration of Statistics, the ISI Centenary Volume*. (eds. Anthony C. Aitkinson and S.E. Fienberg), Springer-Verlag, 543-569.
- Rao, C.R. (1965). On Discrete Distributions Arising out of Methods of Ascertainment. In *Classical and Contagious Discrete Distributions*, (ed. G.P. Patil), Pergamon Press and Statistical Publishing Society, Calcutta, 320-332.
- Rao, C.R. (1973). *Linear Statistical Inference and its Applications*. 2nd ed. Wiley, New York.
- Schucany, W.R., Gray, H.L. and Owen, D.B. (1971). Bias reduction in estimation. *J. Amer. Statist.* 66 524-533.
- Shah, K.R. and Sinha, B.K. (1989). *Theory of Optimal Designs. Lecture notes in Statistics*, 54 (eds. J. Berger, S. Fienberg, J. Gani, K. Krickeberg, and B. Singer ). Springer- Verlag.
- Sen, P.K. (1987). What do the Arithmetic, Geometric and Harmonic means tell us in Length-Biased Sampling. *Statistics & Probability Letters*, 5 95-98.
- Sen, P.K. and Singer, J.M. (1993). *Large Sample Methods in Statistics. An Introduction with Applications*, Chapman & Hall, London.
- Sen, P.K. (1977). Some Invariance Principles to Jackknifing and their Role in Sequential Analysis. *Ann. Statist.* 2 316-329.
- Silvey, S.D. (1980). *Optimal Design. An introduction to the theory for parameter estimation*. University Press, Cambridge.

Zelen, M. (1974). Problems in Cell Kinetics and Early Detection of Disease. *Reliability and Biometry*, SIAM, Philadelphia, 701-726.

Wilks, S.S., (1962). *Mathematical Statistics*. New York: John Wiley.