

BAYES AND LEAST-SQUARES PROCEDURES  
IN SAMPLING FROM FINITE POPULATIONS

by

PRETTA LATHA MEHROTRA

L

Thesis submitted in fulfilment of the requirements  
for the degree of Master of Science

April 1984

To  
Paps & Mi

ACKNOWLEDGEMENTS

The author wishes to express her sincere appreciation and thanks for the valuable guidance and supervision received from Dr. RM. Sekkappan during the course of her research. She also thanks Dr. Ng Vee Ming, her co-supervisor for his efforts and Puan Alauviah Haji Ismail, for typing out this thesis.

## TABLE OF CONTENTS

	<u>PAGE</u>
DEDICATION	ii
ACKNOWLEDGEMENTS	iii
LIST OF NOTATIONS	vi
ABSTRAK	ix
ABSTRACT	xi
CHAPTER 1 - INTRODUCTION	1
1.1 Preliminaries	1
1.2 Sample Survey Model	3
1.3 Likelihood Function	4
1.4 Literature Review	4
1.5 Summary of Results	8
CHAPTER 2 - LEAST-SQUARES PREDICTION APPROACH	9
2.1 Summary	9
2.2 Introduction	9
2.3 Multiple Regression Estimate	10
2.4 Special Cases	14
2.5 Prediction Using Auxiliary Information in Stratification: Multiple Analysis of Covariance Model	28
2.6 Combined Regression Estimate	29
2.7 Special Cases	32
2.8 Separate Regression Estimate	41
2.9 Special Cases	42



	<u>PAGE</u>
CHAPTER 3 - BAYESIAN APPROACH: GENERAL REGRESSION MODEL	50
3.1 Summary	50
3.2 Bayesian Analysis	50
3.3 Regression Model: $p$ auxiliary variables	52
3.4 Special Cases	57
CHAPTER 4 - BAYESIAN APPROACH: STRATIFICATION	64
4.1 Summary	64
4.2 Combined Regression Estimate	64
4.3 Special Cases	66
4.4 Separate Regression Estimate	74
4.5 Special Cases	76
CHAPTER 5 - BAYESIAN APPROACH: DOUBLE SAMPLING FOR STRATIFICATION	85
5.1 Summary	85
5.2 Combined Regression Estimate: $\pi_i$ unknown	85
5.3 Special Cases	91
5.4 Non Response	94
5.5 Combined Regression Estimate Using Two-Phase Sampling: $\pi_i$ unknown	97
APPENDIX I	Summary of Estimates Obtained Under the Least-Squares Prediction Approach
APPENDIX II	Summary of Estimates Obtained Under the Bayesian Prediction Approach
BIBLIOGRAPHY	

# List of Notations

$$U = \{1, 2, \dots, N\}$$

finite population

$$n, N$$

sample and population size respectively

$$\tilde{Y} = (Y_1, \dots, Y_N)$$

response variable vector

$$\bar{Y} = \frac{1}{N} \sum_{k=1}^N Y_k$$

population mean of response variable

$$\bar{X}_{j\cdot} = \frac{1}{N} \sum_{k=1}^N X_{jk}, j = 1, 2, \dots, p$$

population mean of j-th auxiliary variable

$$\bar{X}_{ij\cdot} = \frac{1}{N_i} \sum_{k=1}^{N_i} X_{ijk}, i = 1, 2, \dots, H$$

i-th stratum population mean of j-th auxiliary variable

$$S_y^2 = \sum_{k=1}^N (Y_k - \bar{Y})^2$$

population sum of squares of Y

$$S_{x_j}^2 = \sum_{k=1}^N (X_{jk} - \bar{X}_{j\cdot})^2, j = 1, 2, \dots, p$$

population sum of squares of  $X_{j\cdot}$ .

$$S_{x_j y} = \sum_{k=1}^N (X_{jk} - \bar{X}_{j\cdot})(Y_k - \bar{Y}),$$

population sum of the product of Y and  $X_{j\cdot}$ .

$$j = 1, 2, \dots, p$$

$$R_{xy} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{S_{xy}}{S_x S_y}$$

population correlation coefficient

$$s = (1, \dots, n)$$

sequence of distinct units in sample, s

$$\begin{aligned} s &= (y_1, \dots, y_n) \\ &= (y : y \in s) \end{aligned}$$

sequence of sampled y values

$$S$$

the set of all possible subsets of U

$$s_y^2 = \sum_{k=1}^n (y_k - \bar{y})^2$$

sample sum of squares of y

$$s_{x_j}^2 = \sum_{k=1}^n (x_{jk} - \bar{x}_{j.})^2, j = 1, 2, \dots, p$$

sample sum of squares of  $x_j$ .

$$s_{x_j y} = \sum_{k=1}^n (x_{jk} - \bar{x}_{j.})(y_k - \bar{y}),$$

$j = 1, 2, \dots, p$

sample sum of product of y and  $x_j$ .

$$\rho_{xy} = \frac{s_{xy}}{s_x s_y}$$

sample correlation coefficient of x and y

$$\bar{x}_{ij} = \frac{\sum_{k=1}^{n_i} x_{ijk}}{n_i}, i = 1, 2, \dots, H$$

$j = 1, 2, \dots, p$

i-th stratum sample mean of j-th auxiliary variate

$$\hat{\bar{Y}}$$

$\xi$ -best linear unbiased predictor of mean of Y

$$\hat{\beta}_j, j = 1, 2, \dots, p$$

estimate of regression coefficient

$$\hat{\mu}_i, i = 1, 2, \dots, H$$

estimate of the effect of i-th stratum

$$T_s^t = (1 \dots 1)_{1 \times n}$$

$$T_{s_i}^t = (1 \dots 1)_{1 \times n_i}, i = 1, 2, \dots, H$$

$$T_{\bar{s}}^t = (1 \dots 1)_{1 \times (N-n)}$$

$$T_{\bar{s}_i}^t = (1 \dots 1)_{1 \times (N_i - n_i)}, i = 1, 2, \dots, H$$

$$T_{s_t}^t = (T_{s_1}^t \mid \dots \mid T_{s_H}^t)_{1 \times \sum_{i=1}^H n_i}$$

$$T_{\bar{s}_t}^t = (T_{\bar{s}_1}^t \mid \dots \mid T_{\bar{s}_H}^t) \quad 1 \times \sum_{i=1}^H (N_i - n_i)$$

$$\tilde{m}_{\bar{s}}^t = (N - n, m_{\bar{s}_1}, \dots, m_{\bar{s}_p})$$

$$m_{\bar{s}_j}^t = \sum_{k=1}^{N-n} (x_{jk} - \bar{X}_{j.}), \quad j = 1, 2, \dots, p$$

$$\tilde{m}_{\bar{s}_i}^t = (N_i - n_i, m_{\bar{s}_{i1}}, \dots, m_{\bar{s}_{ip}}), \quad i = 1, 2, \dots, H$$

$$m_{\bar{s}_{ij}} = \sum_{k=1}^{N_i - n_i} (x_{ijk} - \bar{X}_{.j.}), \quad j = 1, 2, \dots, p$$

$$\tilde{m}_{\bar{s}_t}^t = (N_1 - n_1, \dots, N_H - n_H, m_{\bar{s}_{t1}}, \dots, m_{\bar{s}_{tp}})$$

$$m_{\bar{s}_{tj}} = \sum_{i=1}^H \sum_{k=1}^{N_i - n_i} (x_{ijk} - \bar{X}_{.j.}), \quad j = 1, 2, \dots, p$$

## ABSTRAK

Di dalam tesis ini, kami akan pertimbangkan penganggaran dalam populasi-populasi terhingga dengan suatu pendekatan super-populasi berasaskan model, di mana tatacara-tatacara ramalan kuasa dua terkecil dan ramalan Bayesian digunakan. Perhatian yang lebih akan ditumpukan kepada pentaabiran Bayesian. Bab pertama meliputi suatu kajian ringkas am mengenai pendekatan populasi tetap 'klasik' dan pendekatan superpopulasi berasaskan model, kerangka rekabentuk tinjauan am yang digunakan serta suatu tinjauan semula bahan-bahan penulisan yang berkaitan dengan penyelidikan ini.

Di dalam Bab 2, kami kaji semula kesudahan-kesudahan Cassel, Sarndal & Wretman (1977) [Thm. 5.7. Pg. 127] dan Sekkappan (1983), dalam pendekatan ramalan kuasa dua terkecil untuk penerbitan peramal saksama linear terbaik- $\xi$  bagi min populasi,  $\bar{Y}$  dan min ralat kuasa dua terhadap taburan superpopulasi,  $\xi$ . Dengan mengguna kesudahan-kesudahan ini, kami tunjukkan secara berangka, untuk kes-kes tertentu, bagaimana penganggar-penganggar diperolehi. Ditunjukkan juga, secara berangka, bahawa model regresi berganda yang dipertimbangkan akan suai dengan sebarang model regresi polinomial.

Dengan mengikut suatu pendekatan Bayesian dalam Bab 3, kami dapatkan rumus-rumus untuk min dan varians posterior bagi min populasi terhingga dengan mengguna suatu model regresi heteroskedastik linear yang mempunyai  $p$  pembolehubah bantu, dan itlakkan kesudahan-kesudahan Sekkappan (1969 I) untuk prior-prior am atau baur.

Di dalam Bab 4, model regresi linear diperluaskan kepada suatu populasi berstratum. Pendekatan ramalan Bayesian digunakan untuk menganggar min populasi terhingga mengikut suatu model

analisis kovarians berganda dengan ralat heteroskedastik bila terdapat maklumat mengenai lebih dari satu pembolehubah bantu untuk populasi berstratum. Kami pertimbangkan dua kes: bila parameter-parameter regresi yang berkaitan dengan cerapan-cerapan di dalam sebarang stratum adalah

1. sama untuk semua stratum
2. berbeza di antara stratum

Di sini kami tunjukkan bahawa kesudahan-kesudahan dari Sekkappan (1983) boleh diperolehi dari kesudahan-kesudahan Bayesan kami untuk prior baur.

Di dalam model regresi berganda dan model analisis kovarians berganda, sungguhpun anggaran-anggaran yang diperolehi menggunakan pendekatan Bayesan untuk prior-prior baur telah ditunjukkan berhampiran dengan kesudahan-kesudahan bagi pendekatan 'klasik' dan ramalan kuasa dua terkecil, varians bagi anggaran-anggaran akan berbeza dengan pendekatan yang digunakan.

Di dalam Bab 4, adalah dianggapkan bahawa saiz-saiz stratum,  $N_i$  ( $i = 1, 2, \dots, H$ ) diketahui. Untuk keadaan-keadaan tertentu, sungguhpun saiz populasi terhingga,  $N$  diketahui, saiz stratum,  $N_i$  mungkin tidak diketahui. Kesudahan-kesudahan selari yang sepadan dengan penganggar-penganggar yang didapati dalam Bab 4 diterbitkan di dalam Bab 5 melalui pensampelan ganda dua untuk penstratuman. Kesudahan-kesudahan kami mengitlakkan kesudahan-kesudahan Sekkappan (1984) dan Bahadur Singh & Sedransk (1977). Penganggaran min populasi terhingga bila terdapat keadaan tanpa sambutan juga dibincangkan.

## ABSTRACT

In this thesis, we shall consider estimation in finite populations under a model-based superpopulation approach using least-squares and Bayesian prediction procedures. Greater emphasis will be placed on the Bayesian inference. The first chapter covers a general preview of the 'classical' fixed population approach and the superpopulation model-based approach, the general survey design framework used and a literature review relevant to this study.

In Chapter 2, we review, respectively, the result of Cassel, Sarndal & Wretman (1977) [Thm. 5.7, Pg. 127] and Sekkappan (1983) under the least-squares prediction approach towards the derivation of the  $\xi$ -best linear unbiased predictor of the population mean,  $\bar{Y}$  and its mean square error, with respect to the superpopulation distribution,  $\xi$ . Using these results, we shown numerically, for particular cases, how these estimates are obtained. It is also numerically shown that the multiple regression model considered will fit any polynomial regression model.

Follwing a Bayesian approach in Chapter 3, we obtain the formulae for the posterior mean and variance of the finite population mean using a linear heteroscedastic regression model with  $p$  auxiliary variables and hence generalize the results of Sekkappan (1982), Cassel, Sarndal & Wretman (1977) and Ericson (1969 I) for either general or diffuse priors.

In Chapter 4, the linear regression model is extended to a stratified population. Here, the Bayesian prediction approach is used to estimate the finite population mean under a multiple analysis of covariance model with heteroscedastic errors, when information on more than one auxiliary variable is available for stratified populations. We consider two cases: when the regression parameters

associated with the observations in any stratum are

1. the same for all strata
2. different for different stratum.

Here, we show that the results of Sekkappan (1983) can be obtained from our Bayesian results for a diffuse prior.

Under both the multiple regression model and the multiple analysis of covariance model, although the estimates obtained using the Bayesian approach for diffuse priors are shown to agree closely with the usual classical and least-squares prediction results, the mean square errors of these estimates under different approaches will differ.

In Chapter 4, it is assumed that the stratum sizes,  $N_i$ , ( $i=1,2, \dots, H$ ) are known. However, in certain situations, although the finite population size,  $N$  is known, the stratum sizes,  $N_i$  maybe unknown. The parallel results corresponding to the various estimators obtained in Chapter 4 are derived in Chapter 5 using the double sampling technique for stratification. Our results include the results of Sekkappan (1984) and Bahadur Singh & Sedransk (1977) as particular cases. The estimation of the finite population mean when there is non-response is also discussed.



## CHAPTER 1

### INTRODUCTION

#### 1.1 Preliminaries

The estimation procedure in finite population sampling can be studied through two approaches, namely the Fixed Population Approach and the Superpopulation approach. Under the fixed population approach, which is the 'traditional' one in classical survey sampling, with each population unit is associated a fixed but unknown real number, that is, the value of the interested variable. However, following the superpopulation approach, with each population unit is associated a random variable for which a stochastic structure is specified and the actual value associated with a population unit is treated as the outcome of this random variable.

In this thesis, we shall deal with estimation in finite population sampling using Bayesian and Least-squares prediction procedures under the superpopulation approach. The prediction procedures under the superpopulation approach allows the superpopulation distribution,  $\xi$  an essential role in inference. A relevant point to remember in statistical analysis under this approach is that superpopulation models should not be equated with only Bayesian models.

Under the Bayesian approach,  $\xi$  which specifies the joint distribution of  $(Y_1, \dots, Y_N)$  is taken as a prior distribution from which the posterior distribution of the unobserved coordinates of  $(Y_1, \dots, Y_N)$ , given the sample,  $s$ , is derived. Alternatively, under the least-squares prediction approach,  $\xi$  maybe assumed to contain unknown parameters which

must first be estimated using tools of classical inference, such as the least-squares method, the method of maximum likelihood, etc.

Definition (i):  $\hat{\bar{Y}}$  is called the  $\xi$ -unbiased predictor of  $\bar{Y}$ , if and only if, for a given superpopulation distribution,  $\xi$ ,

$$E_{\xi}(\hat{\bar{Y}} - \bar{Y}) = 0, \text{ for all } s \in S$$

Definition (ii):  $\hat{\bar{Y}}$  is called a  $p\xi$ -unbiased predictor of  $\bar{Y}$ , if and only if, for a given sampling design,  $p$  and a superpopulation distribution,  $\xi$ ,

$$E_{\xi}E(\hat{\bar{Y}} - \bar{Y}) = 0, \text{ for all } s \in S$$

In our study, we shall deal with the model-based approach, that is, the approach where with  $\xi$  as the essential element for inference, the sample,  $s$  is treated as given and the design,  $p$  that produced  $s$  is of minor interest. Therefore, our main interest will be to choose  $\hat{\bar{Y}}$ , for a given  $s$ , to minimize  $E_{\xi}(\hat{\bar{Y}} - \bar{Y})^2$  and beyond this, averaging with respect to  $p$  is secondary. This implies that, under the model-based approach, the problem of finding a predictor,  $\hat{\bar{Y}}$ , which is good for any sample actually obtained is more important than just a good strategy  $(p, \bar{Y})$ .

Since the predictor  $\hat{\bar{Y}}_s$  of  $\bar{Y}_s$ , the mean of the unobserved coordinates of the response variable, is related to the predictor  $\hat{\bar{Y}}$  of the overall mean  $\bar{Y}$  through

$$\hat{\bar{Y}} = \frac{1}{N} \left[ n\bar{y} + (N - n)\hat{\bar{Y}}_s \right] \dots\dots\dots (1.1.1)$$

then,  $\hat{\bar{Y}}$  is  $\xi$ -unbiased as a predictor of  $\bar{Y}$ , if and only if, for every  $s \in S$ ,  $\hat{\bar{Y}}_s$  is  $\xi$ -unbiased as a predictor of  $\bar{Y}_s$ .

A more detailed theory on the least squares prediction approach and the Bayesian approach is given in Chapter 2, Section 2.2, and Chapter 3, Section 3.2, respectively.

The Bayesian inference under quadratic loss consists simply of deriving the posterior expectation of  $\bar{Y}$ , given the data  $(s, \underline{y})$ . Hence, we arrive at the Bayes' estimate

$$E[\bar{Y} \mid (s, \underline{y})] = \frac{\sum_{i \in s} y_i + \sum_{i \in \bar{s}} E(Y_i \mid (s, \underline{y}))}{N} \quad (1.1.2)$$

## 1.2 Sample Survey Model

The basic model for survey sampling in the problem of estimation is presented here. A given population,  $U$  consists of  $N$  units indexed by  $k = 1, 2, \dots, N$ . We denote by  $Y_k$  the real variate associated with the unit  $k$ , ( $k = 1, 2, \dots, N$ ) of the population. Then  $\underline{Y} = (Y_1, \dots, Y_N)$  is a point on the Euclidean space,  $R_N$ . Any function of  $R_N$  is called a parametric function, example, the parametric function  $\bar{Y}$  is said to be the population mean if

$$\bar{Y} = \frac{1}{N} \sum_{k=1}^N Y_k, \text{ for all } \underline{Y} \in R_N$$

Our interest will be to estimate  $\bar{Y}$  on the basis of the units sampled from the population  $U$  and the observed values  $y_k$   $k = 1, 2, \dots, n$ , associated with them, that is, on the basis of  $(s, \underline{y}_k : k \in s)$  where  $s$  is a subset of  $U$ , drawn from  $U$  with a given sampling design,  $p$ .

A sampling design,  $p(s)$  is any function of  $S$ , the set of all possible subsets  $s$  of  $U$ , such that  $p(s) \geq 0$  and  $\sum p(s) = 1$ , for  $s \in S$ .

### 1.3 Likelihood Function

Given the sampling design,  $(S, p)$  and the population vector,  $\underline{Y} = (Y_1, \dots, Y_N)$ , the probability of obtaining the data  $(s, y_k : k \in s)$  can be shown as

$$\text{Prob}[(s, y_k : k \in s) | \underline{Y}, (S, p)] = \begin{cases} p(s), & \text{if } Y_k = y_k, \text{ for} \\ & \text{all } k \in s \\ 0, & \text{elsewhere} \end{cases}$$

.....(1.3.1)

Now, with the parameter space,  $R_N$  and parameter vector  $\underline{Y}$ , the likelihood function,  $\ell$  for  $\underline{Y}$  given  $(s, y_k : k \in s)$  is

$$\ell(\underline{Y} | (s, y_k : k \in s)) = \begin{cases} p(s), & \text{if } \underline{Y} \in R_N(y_k : k \in s) \\ 0, & \text{if } \underline{Y} \notin R_N(y_k : k \in s) \end{cases}$$

.....(1.3.2)

where  $\underline{Y} \in R_N(y_k : k \in s)$ , if and only if

$$Y_k = y_k, \text{ for all } k \in s$$

and  $R_N(y_k : k \in s) \subset R_N$

### 1.4 Literature Review

The concept of generalizing from a small 'section' of the population to the 'whole' has been used rather subjectively in various population studies. In order to ascertain how objective the generalization is, two questions which arise are:

- (i) how to select the 'section'
- (ii) how to form the generalization

Basically, the answers to these questions would be to find a combination of selection and estimation techniques which minimize the cost incurred while at the same time optimize the precision of the inference.

During the 1930's, simple ratio and regression estimation were introduced and Cochran (1942) consolidated the theories on this estimation. If we make inferences about a response variable, ignoring the availability of other correlated auxiliary variables, we will sacrifice a lot of information and so Olkin (1958) extended the ratio estimation for finite populations to a multivariate model where he assumed the population means of the auxiliary variables to be known and nonzero. Further, he considered a stratified population for a univariate case and dealt with the two particular ratio estimators, namely

- (i) the separate ratio estimate
- (ii) the combined ratio estimate.

Cochran (1963) compiled the results of Olkin (1958) as well as the estimation under some linear regression models following the fixed population approach. He obtained the combined and separate ratio and regression estimates in a stratified population and adopted the double sampling technique when the population mean of the auxiliary variable is unknown.

When the population means of the auxiliary variables are known, Shukla (1965) derived the multivariate regression estimate. For the case when the population means of auxiliary variables are unknown, Shukla (1965) obtained a double sampling regression estimate.

The inference problem began to take an approach where the finite population was assumed to be generated from an infinite superpopulation. Through prior knowledge, if this superpopulation was known to be of a specific parametric form, example a gamma or normal distribution, then the values of the units not found in the selected sample could be predicted through this distribution of the superpopulation. This is known as predictive statistical inference. Ericson (1969 I & 1969 II) and Scott & Smith (1969) have used this idea in a Bayesian framework for inference.

A non-Bayesian prediction approach utilizing the least-squares method, under some multiple linear regression models in finite populations were considered by Royall (1970), (1971) and Brewer (1963). Cassel, Sarndal and Wretman (1977) derived a  $\xi$ -best linear unbiased predictor of the population mean under a multiple regression model using a model-based approach where  $\xi$  is the superpopulation distribution. Holt & Smith (1979) used the method of least-squares for prediction in a one way analysis of variance model. Under a multiple analysis of covariance model, Sekkappan (1983) generalised the results of Royall (1970), Brewer (1963) and Holt & Smith (1979) to obtain the  $\xi$ -best linear unbiased predictors of  $\bar{Y}$  and their respective mean square errors for the combined and separates regression estimates.

Under the Bayesian approach, Ericson (1969 I) considered estimation under a linear regression model in finite populations. Ericson (1969 I) also gave properties of the posterior distributions for some special cases obtained by assuming the conditional density of a random variate, that is,  $f(y_k | \theta)$ ,

$k = 1, 2, \dots, N$  where  $\theta$  is a real or vector valued parameter, to be a normal with either known or unknown variance and by taking  $dF(\theta)$  to be either a natural conjugate or a uniform prior. In the same paper, Ericson extended his results to the case involving the use of concomitant measurements and showed that under this Bayesian model, several ratio and regression estimators arise as the means of posterior distributions. Later, in Ericson (1969 II), the Bayesian superpopulation approach was extended for stratification. Sekkappan (1981 I) generalized the Bayesian models considered in Ericson (1969 I) for  $K$  characteristics, which was extended to a stratified population in Sekkappan (1981 II). For an analysis of covariance model, Sekkappan (1982) obtained a regression estimate of the population mean when the regression parameters are different in different strata.

Various attempts have been made to identify methods of obtaining a high degree of response since it is evident that the problem of nonresponse exists in most surveys. Bahadur Singh and Sedransk (1977) considered nonresponse in the context of a two phase sampling design, assuming post stratification. This result was generalized by Sekkappan (1984) under a multiple analysis of covariance model when the regression parameters are different for different strata. Earlier, Draper & Guttman (1968 I) and Mohd. Zubair Khan (1976) showed that for a stratified population, a two phase Bayesian sampling scheme will improve the Bayes' estimates of unknown parameters in a regression model.

## 1.5 Summary of Results

This thesis is primarily an extension of the work of Bahadur Singh & Sedransk (1977) and Sekkappan (1981 I), (1982), (1984) for finding the Bayes' estimate for finite population mean. The work also involves an extension of Royall (1970, 1971), in the sense that we have studied some aspects of the estimation problem with numerical examples using the least-squares prediction approach. Each chapter has its own summary at the beginning.



## CHAPTER 2

### LEAST-SQUARES PREDICTION APPROACH

#### 2.1 Summary

In this chapter, we will review the results of Cassel, Sarndal & Wretman (1977) [Thm. 5.7, Pg. 127] and Sekkappan (1983) under the multiple regression model (2.3.1) and the multiple analysis of covariance models (2.6.1) and (2.8.1), respectively, regarding the derivation of the  $\xi$ -best linear unbiased ( $\xi$ -BLU) predictor of  $\bar{Y}$  and its mean square error (MSE) with respect to the superpopulation distribution,  $\xi$ . Based on the results of Cassel, Sarndal & Wretman (1977), the details towards obtaining the  $\xi$ -BLU predictor of  $\bar{Y}$  when  $p = 2$  and  $p = 1$  together with their respective MSE are given. The classical ratio estimator will be deduced by assuming certain constraints in the model (2.3.1). As a special case, the model (2.3.1) is shown to fit any polynomial regression model. For the particular cases of  $p = 2$  and  $p = 1$ , we provide numerical derivations of the  $\xi$ -BLU predictor of  $\bar{Y}$  and its MSE, based on the results of Cassel, Sarndal & Wretman (1977) and Sekkappan (1983).

#### 2.2 Introduction

In the superpopulation approach, the vector of population values  $\underline{y} = (y_1, \dots, y_N)$  is assumed to be the realised outcome of a vector of random variables  $\underline{Y} = (Y_1, \dots, Y_N)$ . The joint distribution of  $Y_1, \dots, Y_N$  will be denoted by  $\xi$ , about which certain features are assumed known. For a superpopulation model, a specified set of conditions define a class of distributions to which the superpopulation distribution,  $\xi$  is

assumed to belong. In the model-based approach, we distinguish two main approaches to inference under superpopulation assumptions:

1. models leading to the use of classical inference tools which we deal with in this chapter
2. models leading to the use of Bayesian inference tools which we deal with in the Chapters 3,4 and 5.

In the classical inference version, the observed  $y_k$  are first used to make inference about a possibly unknown parameter vector  $\theta$  indexing  $\xi$ . Finally,  $\xi$  is used to predict the mean  $\bar{Y}_s$  of the unobserved units.

In this chapter, we shall consider the classical inference of the prediction approach to a multiple regression model and a multiple analysis of covariance model.

### 2.3 Multiple Regression Estimate

The class of probability measures,  $\xi$ , on  $R_N$  is such that  $Y_1, \dots, Y_N$  are independently distributed and

$$E(Y_k) = \mu + \beta_1(X_{1k} - \bar{X}_{1.}) + \dots + \beta_p(X_{pk} - \bar{X}_{p.})$$

$$V(Y_k) = \lambda_k \sigma^2$$

where

$\beta_1, \dots, \beta_p$  and  $\sigma^2$  are unknown

$X_{1k}, \dots, X_{pk}, \lambda_k$  is a set of known numbers for

every  $k(k = 1, 2, \dots, N)$ .

Then, the model for a finite population would be

$$\tilde{y}_{Nx1} = Z_{Nx(p+1)} \beta_{(p+1) \times 1} + \tilde{e}_{Nx1} \quad \dots\dots\dots (2.3.1)$$

where

$$\tilde{y} = (Y_1, \dots, Y_N)^t$$

$$Z = \begin{pmatrix} 1 & (X_{11} - \bar{X}_{1.}) & \dots & (X_{p1} - X_{p.}) \\ \vdots & \vdots & & \vdots \\ 1 & (X_{1N} - \bar{X}_{1.}) & \dots & (X_{pN} - X_{p.}) \end{pmatrix}$$

$$\tilde{\beta} = (\mu, \beta_1, \dots, \beta_p)^t$$

The vector of random disturbances,  $\tilde{e}$ , is distributed with a vector of zero means and a variance covariance matrix  $\sigma^2 \eta_{NxN}^{-1}$  where,

$$\eta_{NxN}^{-1} = \begin{pmatrix} \lambda_1 & . & . & . & 0 \\ \vdots & & & & \vdots \\ 0 & . & . & . & \lambda_N \end{pmatrix} \quad \dots\dots\dots (2.3.2)$$

with  $\lambda_k = g(x_{1k}, \dots, x_{pk})$ ,  $k = 1, 2, \dots, N$

Given a sample of size  $n$ , the model maybe partitioned as

$$\begin{pmatrix} \tilde{y}_{s \times 1} \\ \tilde{y}_{N-n \times 1} \end{pmatrix} = \begin{pmatrix} Z_{s \times p+1} \\ Z_{N-n \times p+1} \end{pmatrix} \tilde{\beta} + \begin{pmatrix} \tilde{e}_{s \times 1} \\ \tilde{e}_{N-n \times 1} \end{pmatrix} \quad \dots\dots\dots (2.3.3)$$

$$\text{with cov } (\tilde{e}) = \sigma^2 \begin{pmatrix} \eta_{n \times n}^{-1} & | & 0_{n \times N-n} \\ \hline 0_{N-n \times n} & | & \eta_{N-n \times N-n}^{-1} \end{pmatrix}$$

Hence, after selecting the sample,  $s$ , of size  $n$ , the model for the observed values is

$$\underset{\sim}{Y}_s = Z_{s\sim} \beta + \underset{\sim}{e}_s \quad \dots\dots\dots (2.3.4)$$

where

$$E(\underset{\sim}{e}_s) = 0$$

$$D(\underset{\sim}{e}_s) = \sigma^2 \eta_s^{-1}$$

By the least-square method, the generalised least-squares estimate of  $\beta$  is

$$\underset{\sim}{\hat{\beta}} = (Z_s^t \eta_s Z_s)^{-1} (Z_s^t \eta_s \underset{\sim}{Y}_s) \quad \dots\dots\dots (2.3.5)$$

and  $D(\underset{\sim}{\hat{\beta}}) = (Z_s^t \eta_s Z_s)^{-1} \sigma^2 \quad \dots\dots\dots (2.3.6)$

$$\hat{\sigma}^2 = \frac{1}{n-(p+1)} \left[ \underset{\sim}{Y}_s^t \eta_s \underset{\sim}{Y}_s - \underset{\sim}{\hat{\beta}}^t Z_s^t \eta_s \underset{\sim}{Y}_s \right] \dots (2.3.7)$$

Under the multiple regression model (2.3.1), the  $\xi$ -BLU predictor of the population mean  $\bar{Y}$ , for any design  $p$ , is given by

$$\begin{aligned} \hat{\bar{Y}} &= \frac{1}{N} \left[ \sum_{k \in s} y_k + \sum_{k \in \bar{s}} \hat{Y}_k \right] \\ &= \frac{1}{N} \left[ n\bar{y} + T_{\bar{s}}^t Z_{\bar{s}} \underset{\sim}{\hat{\beta}} \right] \quad \dots\dots\dots (2.3.8) \end{aligned}$$

where  $\hat{\bar{Y}}_{\bar{s}} = Z_{\bar{s}} \underset{\sim}{\hat{\beta}}$

and  $T_{\bar{s}}^t = (1 \dots 1)_{1 \times N-n}$

We can simplify (2.3.8) further to obtain

$$\hat{\bar{Y}} = \frac{1}{N} \left[ n\bar{y} + (N-n) \underset{\sim}{m}_{\bar{s}}^t \underset{\sim}{\hat{\beta}} \right] \quad \dots\dots\dots (2.3.9)$$

since

$$\begin{aligned} T_{\bar{s}}^t Z_{\bar{s}} &= \left( N-n, \sum_{k=1}^{N-n} (x_{1k} - \bar{X}_{1.}), \dots, \sum_{k=1}^{N-n} (x_{pk} - \bar{X}_{p.}) \right) \\ &= (N-n, m_{\bar{s}1}, \dots, m_{\bar{s}p}) \\ &= \tilde{m}_{\bar{s}}^t \end{aligned}$$

with

$$\begin{aligned} m_{\bar{s}j} &= \sum_{k=1}^{N-n} (x_{jk} - \bar{X}_{j.}) \\ &= \sum_{k=1}^{N-n} x_{jk} - (N-n)\bar{X}_{j.} \\ &= N\bar{X}_{j.} - n\bar{x}_{j.} - (N-n)\bar{X}_{j.} \\ &= n(\bar{X}_{j.} - \bar{x}_{j.}) \quad \text{for } j = 1, 2, \dots, p \end{aligned}$$

The mean square error of  $\hat{\bar{Y}}$ , with respect to the superpopulation distribution  $\xi$ , is

$$\begin{aligned} \text{MSE}(\hat{\bar{Y}}) &= E_{\xi} [\hat{\bar{Y}} - \bar{Y}]^2 \\ &= \frac{1}{N^2} \left[ T_{\bar{s}}^t D(\hat{\bar{Y}}_{\bar{s}}) T_{\bar{s}} \right] \\ &= \frac{1}{N^2} \left[ T_{\bar{s}}^t E \{ (\hat{\bar{Y}}_{\bar{s}} - \bar{Y}_{\bar{s}}) (\hat{\bar{Y}}_{\bar{s}} - \bar{Y}_{\bar{s}})^t \} T_{\bar{s}} \right] \\ &= \frac{1}{N^2} \left[ T_{\bar{s}}^t E \{ (Z_{\bar{s}} \hat{\beta} - Z_{\bar{s}} \beta - \tilde{e}_{\bar{s}}) (Z_{\bar{s}} \hat{\beta} - Z_{\bar{s}} \beta - \tilde{e}_{\bar{s}})^t \} T_{\bar{s}} \right] \\ &= \left[ \frac{1}{N^2} \left[ T_{\bar{s}}^t E \{ Z_{\bar{s}} (\hat{\beta} - \beta) (\hat{\beta} - \beta)^t Z_{\bar{s}}^t + \tilde{e}_{\bar{s}} \tilde{e}_{\bar{s}}^t \} T_{\bar{s}} \right] \right] \\ &= \frac{1}{N^2} \left[ T_{\bar{s}}^t \{ Z_{\bar{s}} \text{cov}(\hat{\beta}) Z_{\bar{s}}^t + \text{cov}(\tilde{e}_{\bar{s}}) \} T_{\bar{s}} \right] \end{aligned}$$

$$\begin{aligned}
 &= \frac{\sigma^2}{N^2} \left[ T_{\bar{s}}^t Z_{\bar{s}} (Z_s^t \eta_s Z_s)^{-1} Z_{\bar{s}}^t T_{\bar{s}} + T_{\bar{s}}^t \eta_{\bar{s}}^{-1} T_{\bar{s}} \right] \\
 &= \frac{\sigma^2}{N^2} \left[ \tilde{m}_{\bar{s}}^t (Z_s^t \eta_s Z_s)^{-1} \tilde{m}_{\bar{s}} + \sum_{k=1}^{N-n} \lambda_k \right] \dots\dots\dots (2.3.10)
 \end{aligned}$$

## 2.4 Special Cases

Case I: Suppose  $\lambda_1 = \lambda_2 = \dots = \lambda_N = 1$ , then  $\eta^{-1}$  becomes an identity matrix. Hence, in model (2.3.1)

$$D(\underset{\sim}{e}) = \sigma^2 I_{N \times N}$$

The estimates (2.3.9) and (2.3.10) will still hold true but

$$\underset{\sim}{\hat{\beta}} = (Z_s^t Z_s)^{-1} Z_s^t \underset{\sim}{y}_s \dots\dots\dots (2.4.1)$$

and

$$\hat{\sigma}^2 = \frac{1}{n-(p+1)} \left[ \underset{\sim}{y}_s^t \underset{\sim}{y}_s - \underset{\sim}{\hat{\beta}}^t Z_s^t \underset{\sim}{y}_s \right]$$

Case II: Suppose  $p = 2$  in Case I, then

$$Z = \begin{pmatrix} 1 & (x_{11} - \bar{x}_{1.}) & (x_{21} - \bar{x}_{2.}) \\ \vdots & \vdots & \vdots \\ 1 & (x_{1N} - \bar{x}_{1.}) & (x_{2N} - \bar{x}_{2.}) \end{pmatrix}$$

$$\underset{\sim}{\beta} = (\mu, \beta_1, \beta_2)^t$$

Given the sample and observed variate values  $y_k (k = 1, 2, \dots, n)$ , the estimates of  $\mu, \beta_1, \beta_2$  can be easily found through (2.4.1) where

$$\begin{aligned}
 (Z_s^t \ Z_s)^{-1} &= \frac{1}{|Z_s^t \ Z_s|} \begin{pmatrix} K_{11} & K_{12} & K_{13} \\ K_{21} & K_{22} & K_{23} \\ K_{31} & K_{32} & K_{33} \end{pmatrix} \\
 &= \begin{pmatrix} K_{11}' & K_{12}' & K_{13}' \\ K_{21}' & K_{22}' & K_{23}' \\ K_{31}' & K_{32}' & K_{33}' \end{pmatrix}
 \end{aligned}$$

with the determinant

$$|Z_s^t \ Z_s| = n s_{x_1}^2 s_{x_2}^2 (1 - \rho_{x_1 x_2}^2)$$

and

$$\begin{aligned}
 K_{11}' &= \frac{1}{n} + \frac{(\bar{x}_{2.} - \bar{X}_{2.})^2}{(1 - \rho_{x_1 x_2}^2) s_{x_2}^2} + \frac{(\bar{x}_{1.} - \bar{X}_{1.})^2}{(1 - \rho_{x_1 x_2}^2) s_{x_1}^2} \\
 &\quad - \frac{2 \rho_{x_1 x_2}^2 (\bar{x}_{1.} - \bar{X}_{1.})(\bar{x}_{2.} - \bar{X}_{2.})}{s_{x_1 x_2} (1 - \rho_{x_1 x_2}^2)}
 \end{aligned}$$

$$K_{12}' = K_{21}' = \frac{\rho_{x_1 x_2}^2 (\bar{x}_{2.} - \bar{X}_{2.})}{s_{x_1 x_2} (1 - \rho_{x_1 x_2}^2)} - \frac{(\bar{x}_{1.} - \bar{X}_{1.})}{s_{x_1}^2 (1 - \rho_{x_1 x_2}^2)}$$

$$K_{13}' = K_{31}' = \frac{\rho_{x_1 x_2}^2 (\bar{x}_{1.} - \bar{X}_{1.})}{s_{x_1 x_2} (1 - \rho_{x_1 x_2}^2)} - \frac{(\bar{x}_{2.} - \bar{X}_{2.})}{s_{x_2}^2 (1 - \rho_{x_1 x_2}^2)}$$

$$K_{22}' = \frac{1}{s_{x_1}^2 (1 - \rho_{x_1 x_2}^2)}$$

$$K_{23}' = K_{32}' = \frac{-\rho_{x_1 x_2}^2}{s_{x_1 x_2} (1 - \rho_{x_1 x_2}^2)}$$

$$K'_{33} = \frac{1}{s_{x_2}^2 (1 - \rho_{x_1 x_2}^2)}$$

On substituting and simplyfying (2.4.1), we obtain

$$\hat{\mu} = \bar{y} - \hat{\beta}_1(\bar{x}_{1.} - \bar{x}_{1.}) - \hat{\beta}_2(\bar{x}_{2.} - \bar{x}_{2.}) \dots\dots\dots (2.4.2)$$

$$\hat{\beta}_1 = \frac{s_{x_1 y}}{s_{x_1}^2 (1 - \rho_{x_1 x_2}^2)} - \frac{\rho_{x_1 x_2}^2 s_{x_2 y}}{s_{x_1 x_2} (1 - \rho_{x_1 x_2}^2)} \dots\dots\dots (2.4.3)$$

$$\hat{\beta}_2 = \frac{s_{x_2 y}}{s_{x_2}^2 (1 - \rho_{x_1 x_2}^2)} - \frac{\rho_{x_1 x_2}^2 s_{x_1 y}}{s_{x_1 x_2} (1 - \rho_{x_1 x_2}^2)} \dots\dots\dots (2.4.4)$$

noting the fact,

$$Z_s^t \tilde{y}_s = \left( \sum_{k \in s} y_k, s_{x_1 y} + (\bar{x}_{1.} - \bar{x}_{1.}) \sum_{k \in s} y_k, s_{x_2 y} + (\bar{x}_{2.} - \bar{x}_{2.}) \sum_{k \in s} y_k \right)^t$$

and  $s_{x_1}^2$ ,  $s_{x_2}^2$ ,  $s_{x_1 y}$ ,  $s_{x_2 y}$ ,  $s_{x_1 x_2}$  and  $\rho_{x_1 x_2}^2$  are defined in the list of notations.

Then, the estimate (2.3.9) becomes

$$\hat{\bar{Y}} = \frac{1}{N} \left\{ n\bar{y} + (N-n) \left[ \hat{\mu} + \frac{\hat{\beta}_1 n(\bar{x}_{1.} - \bar{x}_{1.})}{N-n} + \frac{\hat{\beta}_2 n(\bar{x}_{2.} - \bar{x}_{2.})}{N-n} \right] \right\}$$

which after substituting (2.4.2), (2.4.3) and (2.4.4) will become

$$\hat{\bar{Y}} = \bar{y} + \hat{\beta}_1(\bar{x}_{1.} - \bar{x}_{1.}) + \hat{\beta}_2(\bar{x}_{2.} - \bar{x}_{2.}) \dots\dots\dots (2.4.5)$$

Following this, (2.3.10) becomes

$$MSE(\hat{\bar{Y}}) = \frac{\sigma^2}{N^2} \left[ (N-n) + \mathbf{m}_{\bar{s}}^t (Z_s^t Z_s)^{-1} \mathbf{m}_{\bar{s}} \right] \dots\dots\dots (2.4.6)$$

with  $\mathbf{m}_{\bar{s}}^t = (N-n, n(\bar{x}_{1.} - \bar{x}_{1.}), n(\bar{x}_{2.} - \bar{x}_{2.}))$



After substitution, factorization and simplification, (2.4.6) becomes,

$$\begin{aligned} \text{MSE}(\hat{\bar{Y}}) = \sigma^2 \frac{N-n}{Nn} + \frac{(\bar{x}_{1.} - \bar{X}_{1.})^2}{s_{x_1}^2 (1-\rho_{x_1 x_2}^2)} + \frac{(\bar{x}_{2.} - \bar{X}_{2.})^2}{s_{x_2}^2 (1-\rho_{x_1 x_2}^2)} \\ - \frac{2\rho_{x_1 x_2}^2 (\bar{x}_{1.} - \bar{X}_{1.})(\bar{x}_{2.} - \bar{X}_{2.})}{s_{x_1 x_2}^2 (1-\rho_{x_1 x_2}^2)} \dots\dots\dots (2.4.7) \end{aligned}$$

### Example

In jute fibre crops for estimating mean yield of fibre per plant, two auxiliary variables have been measured, namely, height and base diameter, which are correlated with the yield of fibre. Population consisting of 50 jute plants (Capsulanes) sown at the Jute Agricultural Research Institute Farm, Barrackpore in the year 1962 - 63 were taken and data collected is given in Table 2.1 below. This example is taken from Shukla, G.K. (1965).

Table 2.1

Sr.No.	Plant Hgt. (ft) $x_1$	Base diameter (cm) $x_2$	Fibre weight (g) $y$	Sr.No.	Plant Hgt. (ft) $x_1$	Base diameter (cm) $x_2$	Fibre weight (g) $y$
1	7.08	1.3	5.0	26	6.83	1.5	7.0
2	7.00	1.3	5.5	27	7.17	1.7	7.5
3	7.08	1.4	5.5	28	6.83	1.4	6.5
4	5.67	1.3	3.5	29	7.50	1.8	9.0
5	5.75	1.2	4.0	30	6.25	1.2	5.0
6	6.08	1.3	6.0	31	6.67	1.5	6.0
7	6.83	1.2	4.5	32	6.42	1.4	4.5
8	7.58	1.4	6.5	33	6.25	1.7	4.5
9	6.33	1.5	5.0	34	6.58	1.3	5.0
10	6.17	1.3	4.5	35	7.25	1.2	6.0
11	6.75	1.4	5.5	36	7.33	1.3	6.5
12	6.25	1.3	4.5	37	5.33	1.5	4.5
13	5.92	1.3	3.5	38	6.92	1.4	6.0
14	6.00	1.3	4.5	39	6.75	1.4	5.0
15	7.25	1.5	7.0	40	7.25	1.5	7.0
16	5.50	1.3	4.5	41	5.58	1.6	6.5
17	6.83	1.2	5.5	42	7.42	1.5	7.5
18	6.67	1.2	4.5	43	7.08	1.4	7.0
19	6.83	1.2	5.5	44	7.00	1.7	6.5
20	7.33	1.5	7.5	45	7.33	1.4	7.0
21	5.25	1.4	3.5	46	7.58	1.6	7.5
22	6.17	1.4	4.5	47	6.42	1.7	5.5
23	7.25	1.7	9.5	48	6.83	1.6	7.0
24	7.17	1.6	6.5	49	6.00	1.5	4.0
25	7.08	1.2	5.0	50	6.17	1.5	4.5

$$\bar{x}_{1.} = 6.65 \text{ ft.} \quad \bar{x}_{2.} = 1.42 \text{ cm.} \quad \bar{y} = 5.69 \text{ gm.}$$

To illustrate the method of estimation discussed above, the population units in Table 2.1 are used to select a sample of size of size  $n = 20$ , in order to estimate the mean yield of fibre per plant given information on two auxiliary variables correlated to the response variable. The sample data, given in Table 2.2 below, is selected using the random numbers tables.

Table 2.2

Sr. No.	$x_1$	$x_2$	y	Sr. No.	$x_1$	$x_2$	y
5	5.75	1.2	4.0	29	7.50	1.8	9.0
7	6.83	1.2	4.5	31	6.67	1.5	6.0
9	6.33	1.5	5.0	34	6.58	1.3	5.0
11	6.75	1.4	5.5	35	7.25	1.2	6.0
14	6.00	1.3	4.5	36	7.33	1.3	6.5
16	5.50	1.3	4.5	38	6.92	1.4	6.0
17	6.83	1.2	5.5	41	5.58	1.6	6.5
20	7.33	1.5	7.5	42	7.42	1.5	7.5
26	6.83	1.5	7.0	44	7.00	1.7	6.5
27	7.17	1.7	7.5	49	6.00	1.5	4.0

From the data above, we obtain the following sample values

$$\begin{aligned}\bar{x}_{1.} &= 6.68 & \bar{x}_{2.} &= 1.43 & \bar{y} &= 5.925 \\ s_{x_1}^2 &= 7.3779 & s_{x_2}^2 &= 0.622 \\ s_{x_1 y} &= 11.31 & s_{x_2 y} &= 3.195\end{aligned}$$

$$s_{x_1 x_2} = 0.502$$

$$\rho_{x_1 x_2}^2 = \frac{s_{x_1 x_2}^2}{s_{x_1}^2 s_{x_2}^2} = 0.0549$$

The estimates of  $\beta_1$ ,  $\beta_2$  and  $\mu$  are obtained from (2.4.2), (2.4.3) and (2.4.4), respectively, as

$$\hat{\beta}_1 = \frac{11.31}{7.3779(1-0.0549)} - \frac{0.0549(3.195)}{0.502(1-0.0549)} = 1.2523$$

$$\hat{\beta}_2 = \frac{3.195}{0.622(1-0.0549)} - \frac{0.0549(11.31)}{0.502(1-0.0549)} = 4.1263$$

$$\hat{\mu} = 5.925 - 1.2523(0.03) - 4.1263(0.01) = 5.8462$$

Hence, the estimate of the mean, given by (2.4.5) is

$$\begin{aligned}\hat{\bar{Y}} &= 5.925 + 1.2523(6.65 - 6.68) + 4.1263(1.42 - 1.43) \\ &= 5.8462\end{aligned}$$

whose mean square error, given by (2.4.6), is

$$\begin{aligned} \text{MSE}(\hat{\bar{Y}}) &= \sigma^2 \left[ \frac{30}{1000} + \frac{(0.03)^2}{7.3779(1-0.0549)} + \frac{(0.01)^2}{0.622(1-0.0549)} \right. \\ &\quad \left. - \frac{2(0.0549)(0.03)(0.01)}{0.502(1-0.0549)} \right] \\ &= \sigma^2(0.0302) \end{aligned}$$

where

$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{n-3} \left[ \sum_{k=1}^n y_k^2 - \hat{\mu} \sum_{k=1}^n y_k - \hat{\beta}_1 \sum_{k=1}^n y_k (x_{1k} - \bar{x}_1) - \hat{\beta}_2 \sum_{k=1}^n y_k (x_{2k} - \bar{x}_2) \right] \\ &= \frac{1}{17} \left[ 736.75 - 5.8462(118.5) - 1.2523(14.8650) - 4.1263(4.38) \right] \\ &= 0.4286 \end{aligned}$$

Therefore,

$$\text{MSE}(\hat{\bar{Y}}) = 0.0129$$

Case III: Suppose  $p = 1$  in Case I, then

$$Z = \begin{pmatrix} 1 & \dots & 1 \\ (x_1 - \bar{x}) & \dots & (x_N - \bar{x}) \end{pmatrix}^t$$

$$\beta = (\mu, \beta_1)^t$$

The estimate (2.4.1) will become

$$\begin{aligned} \hat{\beta} &= \begin{pmatrix} n & n(\bar{x} - \bar{X}) \\ n(\bar{x} - \bar{X}) & \sum_{k=1}^n (x_k - \bar{x})^2 \end{pmatrix}^{-1} \begin{pmatrix} \sum_{k=1}^n y_k \\ \sum_{k=1}^n y_k (x_k - \bar{x}) \end{pmatrix} \\ &= (\bar{y} - \hat{\beta}_1(\bar{x} - \bar{X}), \hat{\beta}_1)^t \dots \dots \dots (2.4.8) \end{aligned}$$

after simplification,

$$\text{where } \hat{\beta}_1 = \frac{\sum_{k=1}^n (x_k - \bar{x}) y_k}{\sum_{k=1}^n (x_k - \bar{x})^2}$$

Hence, the estimate (2.3.9) will be

$$\begin{aligned}\hat{\bar{Y}} &= \frac{1}{N} \left[ n\bar{y} + (N-n) \left[ \bar{y} - \hat{\beta}_1 (\bar{x} - \bar{X}) + \frac{\hat{\beta}_1 n (\bar{X} - \bar{x})}{N-n} \right] \right] \\ &= \bar{y} + \hat{\beta}_1 (\bar{X} - \bar{x})\end{aligned}$$

$$\begin{aligned}\text{noting the fact, } Z_s^t \tilde{y}_s &= \left( \sum_{k \in s} y_k, \sum_{k \in s} (x_k - \bar{X}) y_k \right)^t \\ &= \left( \sum_{k \in s} y_k, s_{xy} + (\bar{x} - \bar{X}) \sum_{k \in s} y_k \right)^t\end{aligned}$$

The estimate (2.4.9) is known in classical survey sampling as the regression estimator which is also obtained by Cassel, Sarndal & Wretman (1977), Example 2, Section 5.6.

The mean square error of  $\hat{\bar{Y}}$  (2.4.9) will be

$$MSE(\hat{\bar{Y}}) = \sigma^2 \frac{(N-n)^2}{N^2} \left[ m_{\tilde{s}}^t (Z_s^t Z_s)^{-1} m_{\tilde{s}} + \frac{1}{N-n} \right] \dots\dots\dots (2.4.10)$$

$$\text{where } m_{\tilde{s}}^t = \left( 1, \frac{n(\bar{X} - \bar{x})}{N-n} \right)$$

Formula (2.4.10) can be further simplified to obtain

$$MSE(\hat{\bar{Y}}) = \sigma^2 \left[ \frac{N-n}{Nn} + \frac{(\bar{X} - \bar{x})^2}{2 s_x^2} \right]$$

### Example

We use the same data given in Table 2.2, but for our interest, let the response variate,  $y$  be related to only one auxiliary variable, namely, plant height. Therefore, for a sample of size  $n = 20$ , the data given in Table 2.3 will be used

Table 2.3

Sr.No.	x	y	Sr.No.	x	y
5	5.75	4.0	29	7.50	9.0
7	6.83	4.5	31	6.67	6.0
9	6.33	5.0	34	6.58	5.0
11	6.75	5.5	35	7.25	6.0
14	6.00	4.5	36	7.33	6.5
16	5.50	4.5	38	6.92	6.0
17	6.83	5.5	41	5.58	6.5
20	7.33	7.5	42	7.42	7.5
26	6.83	7.0	44	7.00	6.5
27	7.17	7.5	49	6.00	4.0

We obtain the following sample values

$$\bar{x} = 6.680$$

$$\bar{y} = 5.925$$

$$s_x^2 = 7.3779$$

$$s_{xy} = 11.31$$

From these values, we obtain the estimates of

$$\hat{\beta}_1 = \frac{s_{xy}}{s_x^2} = \frac{11.31}{7.3779} = 1.5330$$

and from (2.4.8)

$$\hat{\mu} = 5.925 - 1.5330(6.68 - 6.65) = 5.8790$$

The estimate of the mean, as given by (2.4.9), is

$$\begin{aligned}\hat{\bar{Y}} &= 5.925 + 1.5330(6.65 - 6.68) \\ &= 5.8790\end{aligned}$$

whose mean square error, as given by (2.4.11), is

$$\begin{aligned}\text{MSE}(\hat{\bar{Y}}) &= \sigma^2 \left[ \frac{30}{1000} + \frac{(0.03)^2}{7.3779} \right] \\ &= \sigma^2 (0.0312)\end{aligned}$$

$$\begin{aligned}
 \text{where } \hat{\sigma}^2 &= \frac{1}{n-2} \left[ \sum_{k=1}^n y_k^2 - \hat{\mu} \sum_{k=1}^n y_k - \hat{\beta}_1 \sum_{k=1}^n y_k (x_k - \bar{x}) \right] \\
 &= \frac{1}{18} [736.75 - 5.8790(118.5) - 1.5530(14.865)] \\
 &= 0.9611
 \end{aligned}$$

Therefore,

$$MSE(\hat{\bar{Y}}) = 0.03$$

Case IV: Suppose we let  $\mu = \beta_1 \bar{X}_1 + \beta_2 \bar{X}_2 + \dots + \beta_p \bar{X}_p$  in model (2.3.1), then the model becomes a regression through the origin, that is

$$\tilde{Y}_{N \times 1} = Z_{N \times p}^* \beta_{p+1}^* + \tilde{e}_{N \times 1} \quad \dots\dots\dots (2.4.12)$$

where

$$\begin{aligned}
 Z^* &= \begin{pmatrix} X_{11} & \dots & X_{p1} \\ \vdots & & \\ X_{1N} & \dots & X_{pN} \end{pmatrix} \\
 \tilde{\beta}^* &= (\beta_1, \dots, \beta_p)^t
 \end{aligned}$$

The results (2.3.9) and (2.3.10) still hold good for model (2.4.12).

Case V: In Case IV, if  $\eta^{-1}$  is an identity matrix, that is,  $\lambda_k = 1$  for  $k = 1, 2, \dots, N$ , the model (2.4.12) becomes a linear homoscedastic model.

The estimate (2.3.9) becomes

$$\hat{\bar{Y}} = \frac{1}{N} \left[ n\bar{y} + (N-n) \tilde{m}_{\bar{s}}^{*t} \tilde{\beta}^* \right] \quad \dots\dots\dots (2.4.13)$$

$$\text{where } \tilde{m}_{\bar{s}}^{*t} = (\tilde{m}_{\bar{s}1}^*, \dots, \tilde{m}_{\bar{s}p}^*)$$

$$\text{with } m_{\bar{s}_j}^* = \frac{N\bar{X}_{j\cdot} - n\bar{x}_{j\cdot}}{N-n}$$

$$\text{and } \hat{\beta}_{\sim}^* = (Z_s^{*t} Z_s^*)^{-1} Z_s^{*t} Y_{\sim s}$$

while (2.3.10) becomes

$$\text{MSE}(\hat{\bar{Y}}) = \frac{\sigma^2}{N^2} \left[ (N-n)^2 m_{\bar{s}}^{*t} (Z_s^{*t} Z_s^*)^{-1} m_{\bar{s}}^* + (N-n) \right] \dots\dots\dots (2.4.14)$$

$$\text{where } \hat{\sigma}^2 = \frac{1}{n-p} \left[ Y_{\sim s}^t Y_{\sim s} - \hat{\beta}_{\sim}^{*t} Z_s^{*t} Y_{\sim s} \right]$$

Case VI: Suppose  $p=1$  in Case IV, and  $\lambda_k = g(x_k) = x_k$ ,  
 $k = 1, 2, \dots, N$ . Then

$$\eta^{*-1} = \begin{pmatrix} x_1 & \dots\dots\dots 0 \\ \vdots & & \vdots \\ 0 & \dots\dots\dots x_N \end{pmatrix}$$

and the estimate (2.3.8) becomes

$$\begin{aligned} \hat{\bar{Y}} &= \frac{1}{N} \left[ n\bar{y} + \sum_{k=1}^{N-n} x_k \cdot \hat{\beta}_1 \right] \\ &= \frac{1}{N} \left[ n\bar{y} + \hat{\beta}_1 (X - n\bar{x}) \right] \dots\dots\dots (2.4.15) \end{aligned}$$

$$\begin{aligned} \text{where } \hat{\beta}_1 &= (Z_s^{*t} \eta_s^* Z_s^*)^{-1} Z_s^{*t} \eta_s^* Y_{\sim s} \\ &= \frac{\bar{y}}{\bar{x}} \dots\dots\dots (2.4.16) \end{aligned}$$

Using (2.4.16) and (2.4.15), we get

$$\hat{\bar{Y}} = \frac{\bar{y}}{\bar{x}} \bar{X} \dots\dots\dots (2.4.17)$$

which is the classical ratio estimator of the  
 population mean.