## ESTIMATION BASED ON CONDITIONAL SPECIFICATION

Grace O. Esimai and Chien-Pai Han Iowa State University

## 1. Introduction

Consider a (p + 1) random vector which follows a multivariate normal distribution where Y is a scalar and X is a p x l vector  $(p \ge 1)$ . In estimating the population mean  $\mu_y$  of Y, it is well known that the precision of the estimator can be increased if auxiliary information is available. In this paper, we shall consider the linear regression estimator of  $\mu_{y}$  with <u>X</u> as the auxiliary variable. To use the regression estimator we need to know the population mean  $\underline{\mu}_{X}$  of  $\underline{X}$ . When  $\underline{\mu}_{\chi}$  is unknown, we may take a preliminary sample to estimate it. This sampling procedure is the double sampling technique. In certain situations, an investigator may have partial information about  $\underline{\mu}_X$  and suspects that  $\underline{\mu}_{X} = \underline{\mu}_{O}$ . In order to utilize this partial information, the investigator can perform a preliminary test about the hypothesis  $H_0: \mu_X = \mu_0$  versus  $H_1: \mu_X \neq \mu_0$ . As an example, consider estimating the average yield per acre of a certain crop. It is known that the yield is highly correlated with the moisture and nitrogen content of the soil. Hence, the moisture and nitrogen content can be used as the auxiliary variable,  $\underline{X}$ . The experimenter usually does not know  $\mu_{\chi}$ ; but from the amount of rainfall reported by the weather bureau or other sources and from analysis by the soil science department, he believes that  $\mu_x$ should be  $\mu_0$ . Once a preliminary sample is available, the investigator may test  $H_0$ . He then will use  $\mu_0$  in the regression estimator if  $H_0$  is accepted; otherwise he uses the sample mean based on the preliminary sample. This estimator is usually known as the preliminary test estimator. If the investigator's prior information or experience is reliable, then the true mean  $\underline{\mu}_{\underline{X}}$  of  $\underline{X}$  will be expected to be very close to  $\mu_0$ . In this situation, the efficiency of the preliminary test estimator is high. Thus in practice, it is desirable to use the preliminary test estimator when partial

Preliminary test estimator was first studied by Bancroft (1944). It belongs to the area of inference based on conditional specification. A bibliography on inference based on conditional specification was recently compiled by Bancroft and Han (1977).

information is available to the investigator.

Let  $\begin{pmatrix} \underline{Y} \\ \underline{X} \end{pmatrix} \sim N(\underline{\mu}, \underline{\Sigma}); \quad \underline{\mu} = \begin{pmatrix} \mu \\ \underline{\mu} \\ \underline{X} \end{pmatrix}$  and  $\underline{\Sigma} = \begin{pmatrix} \sigma^2 & \underline{\Sigma}_{12} \\ \underline{\Sigma}_{21} & \underline{\Sigma}_{22} \end{pmatrix}$ . We assume  $\underline{X}$  is cheaply observed while Y is more expensive to observe. Let  $(Y_1, X_{11}, X_{21}, \dots, X_{p1})'$  i = 1, ...,  $n_2$ be a random sample from  $N(\underline{\mu}, \underline{\Sigma})$ . This is supplemented by  $n_1 - n_2 (n_1 > n_2)$  more independent observations on  $\underline{X} = (X_1, \dots, X_p)'$ . In practice, the sample of  $n_2$  observations is usually a subsample from the sample of  $n_1$ observations. From all the observations, we define  $n_1 \qquad n_1$   $\underline{X}_1 = (1/n_1)(\underline{\Sigma} X_{11}, \dots, \underline{\Sigma} X_{p1})'$ , and from i=1  $n_2$ the subsample, we define  $\overline{y} = 1/n_2 \ \underline{\Sigma} \ y$ , and  $\underline{X}_2 = (\underline{\mu}/n_2)(\underline{\Sigma} X_{11}, \dots, \underline{\Sigma} X_{p1})'$ . If the i=1  $n_2$ and  $\underline{X}_2 = (\underline{\mu}/n_2)(\underline{\Sigma} X_{11}, \dots, \underline{\Sigma} X_{p1})'$ . If the vector  $\underline{\mu}_Y$  and  $\underline{\Sigma}$  are known, then given  $\underline{X}$  an un-

biased estimator of  $\mu_v$  is  $\hat{\mu}_v|_X$ 

$$= \bar{y} + \underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} (\underline{\mu}_{X} - \underline{\bar{X}}_{2}) \text{ with variance}$$

 $(1/n_2) \{ \sigma^2 - \Sigma_{12} \Sigma_{22} \Sigma_{21} \}$ . If  $(1/n_2) \Sigma_{12} \Sigma_{22} \Sigma_{21}$  is considerably large, we have an appreciable gain in precision.

If  $\underline{\mu}_X$  is unknown and partial information about  $\underline{\mu}_X$  is available, without loss of generality we let  $\underline{\mu}_O = \underline{O}$ , the linear regression preliminary test estimator is defined as

$$\hat{\mu}_{lr} = \begin{cases} \bar{y} - \underline{\Sigma}_{12} \ \underline{\Sigma}_{22}^{-1} \ \underline{\bar{x}}_{2} \\ \text{if } n_{1} \ (\underline{\bar{x}}_{1}^{'} \ \underline{\Sigma}_{22}^{-1} \ \underline{\bar{x}}_{1}) \le x_{p,\alpha}^{2} \\ \\ \bar{y} + \underline{\Sigma}_{12} \ \underline{\Sigma}_{22}^{-1} \ (\underline{\bar{x}}_{1} - \underline{\bar{x}}_{2}) \\ \\ \text{if } n_{1} (\underline{\bar{x}}^{'}_{1} \ \underline{\Sigma}_{22}^{-1} \ \underline{\bar{x}}_{1}) > x_{p,\alpha}^{2} \end{cases}$$
(1.1)

where  $\chi^2_{p,\alpha}$  is the 100(1- $\alpha$ ) percent point of the Chi-squared distribution with p degrees of freedom and  $\alpha$  is the level of significance of the preliminary test. Han (1973) studied the estimator  $\hat{\mu}_{lr}$  when p = 1. This paper will consider the general case when  $p \geq 1$ . The bias, mean squared error (MSE) and relative efficiency of  $\stackrel{\wedge}{\mu_{gr}}$  are derived in Esimai (1977) and are given in Section 2. The optimal sample design is discussed in Section 3.

When  $\Sigma$  is unknown, the linear regression preliminary test estimator is

$$\widetilde{\mu}_{\ell r} = \begin{cases} \overline{y} - \underline{s}_{12} \ \underline{s}_{22}^{-1} \ \underline{\bar{x}}_{2} \\ \text{if } m_{1} n_{1} (\underline{\bar{x}}_{1} \underline{s}_{22}^{-1} \ \underline{\bar{x}}_{1}) \leq T_{0}^{2} \\ \overline{y} + \underline{s}_{12} \ \underline{s}_{22}^{-1} (\underline{\bar{x}}_{1} - \underline{\bar{x}}_{2}) \\ \text{if } m_{1} n_{1} (\underline{\bar{x}}_{1} \ \underline{s}_{22}^{-1} \ \underline{\bar{x}}_{1}) > T_{0}^{2} \end{cases}$$
(1.2)

where  $m_1 = n_1 - 1$ ,  $T_0^2$  is the  $100(1-\alpha)$ th percentile of the Hotelling's  $T^2$  distribution with  $m_1$  degrees of freedom. We define

$$\underline{\underline{S}} = \begin{pmatrix} \underline{\underline{S}}_{11} & \underline{\underline{S}}_{12} \\ \underline{\underline{S}}_{21} & \underline{\underline{S}}_{22} \end{pmatrix} \text{ where } \underline{S}_{11} = \sum_{i=1}^{\Sigma} (\underline{y}_i - \overline{y}),$$

$$\underline{\underline{S}}_{12} = \sum_{\substack{i=1\\i=1}}^{n_2} (\underline{y}_i - \overline{y})(\underline{x}_i - \underline{x}_2)',$$

$$\underline{\underline{S}}_{22} = \sum_{\substack{i=1\\i=1}}^{n_1} (\underline{x}_i - \overline{\underline{x}}_1)(\underline{x}_i - \overline{\underline{x}}_1)' \text{ and } \overline{y}, \quad \underline{\overline{x}}_1 \text{ and}$$

$$\underline{\overline{X}}_2 \text{ are as defined above.}$$

2. Bias, MSE and Relative Efficiency of 
$$\hat{\mu}_{lr}$$

The joint distribution of  $(\underline{\bar{x}}_1', \underline{\bar{x}}_2', \bar{y})$ is normal. Denote the acceptance region for the preliminary test by A and its complement by  $\bar{A}$  and let  $\chi^2_{p,\alpha} = b$ .  $E(\hat{\mu}_{gr}) = E \{(\bar{y} - \underline{\Sigma}_{12} \ \underline{\Sigma}_{22}^{-1} \ \underline{\bar{x}}_2) | A\} P(A)$ 

+ E { [
$$\overline{\mathbf{y}} + \underline{\Sigma}_{12} \ \underline{\Sigma}_{22}^{-1} \ (\underline{\overline{\mathbf{x}}}_1 - \underline{\overline{\mathbf{x}}}_2)$$
] ] ] ]   
=  $\mu_{\mathbf{y}} - \underline{\Sigma}_{12} \ \underline{\Sigma}_{22}^{-1} \ \mu_{\mathbf{X}} + \underline{\Sigma}_{12} \ \underline{\Sigma}_{22}^{-1} \ E{\{\underline{\overline{\mathbf{x}}}_1\}}$    
(2.1)

 $B_{l}$  is evaluated in Esimai (1977) and found to be

 $B_{1} = -\underline{\Sigma}_{12} \ \underline{\Sigma}_{22}^{-1} \ \underline{H}_{X} \ H_{p+2} \ (b; \ \delta)$ (2.2)

where  $H_{p+2}(b; \delta)$  is the cumulative distribution function of the noncentral chi-squared distribution with p+2 degrees of freedom and noncentrality parameter  $\delta = n\underline{\mu}'_X \sum_{22}^{-1} \underline{\mu}_X$ .

Without loss of generality, we let  $\Sigma_{22} = I$ and  $\sigma^2 = 1$ . Since  $B_1$  changes sign with  $\Sigma_{12}$ and  $\mu_{\chi}$ , we need only study the bias for  $\mu_{\mathbf{v}} > 0$  and  $\rho > 0$  for p = 1. The bias was also studied by Han (1973) where the bias was expressed in terms of the cumulative distribution function of the standard normal distribution. The two expressions are equivalent as they should be. The general behavior of  $-B_1$ is as follows. The bias is zero when  $\mu_{\chi} = 0$ which is when the null hypothesis is true. It is an increasing function of  $\rho$ , but a decreasing function of  $\alpha$ . For fixed  $n_1$ ,  $\alpha$ and  $\rho$ , the bias increases from zero and then decreases to zero as  $\ \mu_{X}$  increases from zero to one. The values of  $-B_1$  for  $n_1 = 30$ , p = 2and certain values of  $\underline{\Sigma}_{12}$ ,  $\underline{\mu}_{\chi}$  and  $\alpha$  are given in Table 1. The properties of the bias are found to be identical with those recorded for p = 1.

(5 (.7)					
$\alpha = .05$					
0 05 .02 0					
$\alpha$ = .10					
0 03 .01 0					

The MSE of 
$$\hat{\mu}_{\ell r}$$
 was found to be  $M_1 = MSE(\hat{\mu}_{\ell r})$   
=  $g_1 + h_1$  where  
 $g_1 = (1/n_2)\sigma^2 + (1/n_1)\underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} \underline{\Sigma}_{21}$   
 $-(1/n_2)\underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} \underline{\Sigma}_{21}$ , (2.3)  
 $h_1 = \underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} \underline{\mu}_X \underline{\mu}_X' \underline{\Sigma}_{22} \underline{\Sigma}_{21}$   
 $-(1/n_1)\underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} \underline{\Sigma}_{21} H_{p+2}$  (b;  $\delta$ )  
 $-2 \underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} \underline{\mu}_X \underline{\mu}_X' \underline{\Sigma}_{22}^{-1} \underline{\Sigma}_{21} [1 - H_{p+2} (b; \delta)]$   
 $+ \underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} \underline{\mu}_X \underline{\mu}_X' \underline{\Sigma}_{21}^{-1} \underline{\Sigma}_{21} [1 - H_{p+4} (b; \delta)].$ 

Now we compare the performance of the preliminary test estimator,  $\hat{\mu}_{\ell r}$  with the usual linear regression estimator,  $\bar{y} + \underline{\Sigma}_{12} \ \underline{\Sigma}_{22}^{-1} \ (\underline{\bar{X}}_1 - \underline{\bar{X}}_2)$ ,

when the information of  $\mu_X$  is ignored. The relative efficiency of  $\hat{\mu}_{\ell r}$  to the linear regression estimator is defined as

$$e_{1} = \frac{MSE \left(\bar{y} + \underline{\Sigma}_{12} \ \underline{\Sigma}_{22}^{-1} \ (\underline{\bar{x}}_{1} - \underline{\bar{x}}_{2})\right)}{MSE \ (\hat{\mu}_{\ell r})} = \frac{g_{1}}{g_{1} + h_{1}} (2.4)$$

Table 2.	Values of n <sub>2</sub> = 10	e <sub>l</sub> for	p = 2, n	= 30,
<u>Ξ</u> 12	$\begin{pmatrix} 0\\ 0 \end{pmatrix}$	$\left(\begin{array}{c} \cdot 5\\ 0\end{array}\right)$	(•5) (•5)	(1.0) 0)
		α = .05		
(.7, 0) (.5, .5) (.7, .7) ( <b>-</b> .5, .7)	1.24 1.25 4.07 1.64	.69 .83 .56 .80	.70 .52 .22 1.04	1.0 1.0 1.0 1.0
		α = .10		
(.7, 0) (.5, .5) (.7, .7) ( <b>-</b> .5, .7)	1.19 1.20 2.71 1.48	.78 .88 .67 .87	.79 .64 .32 1.02	1.0 1.0 1.0 1.0
		α = .25		
(.7, 0) (.5, .5) (.7, .7) ( <b>-</b> .5, .7)	1.11 1.11 1.61 1.24	.91 .95 .85 .95	.92 .83 .57 1.01	1.0 1.0 1.0 1.0

Without loss of generality we let  $\underline{\Sigma}_{22} = I$  and  $\sigma^2 = 1$ . The values of  $e_1$  for p = 1 are given in Han (1973) and will not be given here. The values of  $e_1$  for p = 2,  $n_1 = 30$ ,  $n_2 = 10$  and certain values of  $\underline{\Sigma}_{12}$ ,  $\alpha$  and  $\underline{\mu}_X$  are given in Table 2. It is seen that  $e_1$  assumes maximum value at  $\underline{\mu}_X = \underline{0}$ . The maximum value of  $e_1$  is an increasing function of  $\rho$  for fixed  $\alpha$ ,  $n_1$  and  $n_2$ . The value of  $e_1$  decreases to a minimum and then increases to unity as  $\underline{\mu}_X'$  increases from (0, 0).

The estimator  $\tilde{\mu}_{\ell r}$  in (1.2) is given when  $\underline{\Sigma}$  is unknown. The bias,  $B_2$ , and the mean square error,  $M_2$ , are derived in Esimai (1977) and are omitted here. The behavior of  $B_2$  is the same as that of  $B_1$  and the behavior of  $M_2$  is similar to that of  $M_1$ .

## 3. The Optimal Sample Design

We shall now consider the problem of finding the optimum allocation of the sample sizes n<sub>1</sub> and n<sub>2</sub> for the estimator  $\hat{\mu}_{lr}$  the cost function is

$$C = n_1 c_1 + n_2 c_2$$
 (3.1)

where  $c_1$  is the cost of observing the vector  $\underline{X}$  and  $c_2$  is the cost of observing Y. The optimum values of  $n_1$  and  $n_2$  are obtained by minimizing  $M_1$  subject to the constraint (3.1). We recall that in practice, under the supposition of a conditional specification, the experimenter has only partial information based on which he believes that  $\underline{\mu}_X$  is close to  $\underline{0}$ . The relative efficiency of  $\hat{\mu}_{\ell r}$  is the largest at  $\underline{\mu}_X = \underline{0}$  and so it would be reasonable to consider the problem of optimum allocation under the optimum situation by letting  $\underline{\mu}_X = \underline{0}$  in  $M_1$ . When  $\underline{\mu}_X = \underline{0}$ ,  $M_1$  becomes

$$M_1 = k_1/n_1 + k_2/n_2$$

(3.2)

$$k_{1} = \underline{\Sigma}_{12} \ \underline{\Sigma}_{22}^{-1} \ \underline{\Sigma}_{21} \ [1 - H_{p+2} \ (b; 0)]$$

$$k_{2} = \sigma^{2} - \underline{\Sigma}_{12} \ \underline{\Sigma}_{22}^{-1} \ \underline{\Sigma}_{21}$$
Minimizing (3.2) subject to (3.1) we find

$$m_{1} = \frac{C\sqrt{k_{1}}}{\sqrt{k_{2}c_{1}c_{2}} + c_{1}\sqrt{k_{1}}},$$

$$m_{2} = \frac{C\sqrt{k_{2}}}{\sqrt{k_{1}c_{1}c_{2}} + c_{2}\sqrt{k_{2}}}$$
(3.3)

and the optimum value of  $M_1$  is

where

$$M_{1, \text{ opt}} = \frac{\left(\sqrt{k_{1}c_{1}} + \sqrt{k_{2}c_{2}}\right)^{2}}{C}$$
(3.4)

We now compare  $M_1$ , opt with the optimum value of the MSE of  $\bar{y} + \underline{\Sigma}_{12} \ \underline{\Sigma}_{22}^{-1} (\underline{\bar{x}}_1 - \underline{\bar{x}}_2)$ , the regression estimator under double sampling without using the preliminary test. If we denote the MSE of  $\bar{y} + \underline{\Sigma}_{12} \ \underline{\Sigma}_{22}^{-1} (\underline{\bar{x}}_1 - \underline{\bar{x}}_2)$  by M

$$M = k_1'/n_1 + k_2'/n_2$$
 (3.5)

where  $k'_1 = \sum_{12} \sum_{22}^{-1} \sum_{21}^{-1}$ ,  $k'_2 = \sigma^2 - \sum_{12} \sum_{22}^{-1} \sum_{21}^{-1}$ and the optimum value of M is

$$M_{opt} = \frac{(\sqrt{k_1'c_1} + \sqrt{k_2'c_2})^2}{C}$$
(3.6)

To compare (3.4) and (3.6) we note from (3.2) that  $(1 - H_{p+2}(b;0))$  is a decreasing function of b with a maximum equal to unity at b = 0. Hence the numerator of  $M_{1, opt}$  at most as large as that of  $M_{opt}$  and  $M_{1, opt} \leq M_{opt}$  with equality holding for b = 0, i.e. when the two estimators coincide. References

- Bancroft, T. A. (1944). On biases in estimation due to the use of preliminary tests of significance. <u>Ann. Math. Stat. 15</u>, 190-204.
- Bancroft, T. A. and Han, C. P. (1977). Inference based on conditional specification: A note and a bibliography. Accepted for publication in <u>International</u> <u>Statistical Review</u>.
- Esimai, Grace O. (1977). Regression estimation for multivariate normal distributions. Unpublished Ph.D. thesis, Iowa State University, Ames, Iowa.
- Han, C. P. (1973). Double sampling with partial information on auxiliary variables. J. Amer. Statist. Assoc. 68, 914-918.