

This PDF is a selection from an out-of-print volume from the National Bureau of Economic Research

Volume Title: Annals of Economic and Social Measurement, Volume 2, number 4

Volume Author/Editor: Sanford V. Berg, editor

Volume Publisher:

Volume URL: <http://www.nber.org/books/aesm73-4>

Publication Date: October 1973

Chapter Title: Criteria Constraints And Multicollinearity In Random Coefficient Regression Models

Chapter Author: P. A. V. B. Swamy

Chapter URL: <http://www.nber.org/chapters/c9935>

Chapter pages in book: (p. 429 - 450)

## CRITERIA, CONSTRAINTS AND MULTICOLLINEARITY IN RANDOM COEFFICIENT REGRESSION MODELS

BY P. A. V. B. SWAMY\*

*This paper analyzes six alternative estimators for random coefficient regression models: (1) minimum variance linear unbiased estimator (MVLU), (2) the Stein-like estimator, (3) the ridge regression estimator, (4) minimum conditional mean square error estimator (MC MSE), (5) the mixed regression estimator, and (6) a maximum likelihood estimator (ML). Attention is focused on the criteria of estimation and parametric constraints in RCR models.*

### I. INTRODUCTION

It has been recognized by many econometricians that the usefulness of the conventional fixed-parameter regression model in the analysis of cross-section data is limited because individuals differ greatly in their behavior, and the diversity of individual decision units implies parameter variation across units, see Swamy (1971) and the references cited therein. In recent years, econometric models, which permit different schemes of parameter variation, have been developed. All these different schemes have been compared by Swamy (1972) who developed an asymptotically efficient procedure of estimating the parameters in a general random coefficient regression (RCR) model. Application of these estimation methods in the analysis of real world data is just beginning, see Feige and Swamy (1972). It has been observed that the use of RCR methods can result in more fruitful and meaningful econometric analyses of micro panel data. In the present paper we analyze alternative estimators with purely algebraic tools. Attention is focused on the criteria of estimation and parametric constraints in RCR models.

The plan of the paper is as follows. Section 2 sets out the estimation rules for random coefficient regression models with and without an unbiasedness condition. Constraints on the parameters and partial prior information are introduced in Section 3 and it is indicated how their presence can help estimation. Methods of using sample data in conjunction with the first two moments of a prior distribution are reviewed in Section 4. The maximum likelihood method of estimating the parameters of a random coefficient model is discussed in Section 5. Summary and Conclusions of the study are presented in Section 6.

### 2. RANDOM COEFFICIENT REGRESSION MODEL

#### 2.1. *The Model*

Swamy (1971) considers the problem of estimating the following equation from a time series of cross-sections.

$$(1) \quad y_i = X_i \beta_i + u_i \quad (i = 1, 2, \dots, n)$$

\* I am grateful to Professor A. Zellner and Dr. Richard D. Porter for helpful comments.

where  $y_t \equiv (y_{t1}, y_{t2}, \dots, y_{tT})'$  is a  $T \times 1$  vector of observations on a dependent variable,  $X_t \equiv (X_{tk})$  ( $k = 1, 2, \dots, K; t = 1, 2, \dots, T$ ) is a  $T \times K$  matrix of observations on  $K$  independent variables,  $\beta_t$  is a  $K \times 1$  vector of coefficients, and  $u_t \equiv (u_{t1}, u_{t2}, \dots, u_{tT})'$  is a  $T \times 1$  vector of disturbances.

Observations on  $y$ 's and  $x$ 's for  $n$  individuals taken over  $T$  periods of time are available. These temporal cross-section data are obtained by assembling cross-sections of  $T$  years, with the same  $n$  cross-section units appearing in all years. The individuals here may be firms, consumers or regions. The subscript  $i$  indexes cross-section observations and the subscript  $t$  indexes time series observations.

In (1) both  $\beta_t$  and  $u_t$  are regarded as realizations of random vectors,<sup>1</sup> and the following assumptions are made.

*Assumption 1:*

(1) The rank of  $X_t$  is  $K$ ,  $n > K$  and  $T > K$ ;

(2) For  $i, j = 1, 2, \dots, n$ :  $E u_i = 0$  and  $E u_i u_j' = \sigma_{ij} \Omega_{ij}$  where

$$\Omega_{ij} = \frac{1}{1 - \rho_i \rho_j} \begin{bmatrix} 1 & \rho_i & \rho_i^2 & \dots & \rho_i^{T-1} \\ \rho_j & 1 & \rho_j & \dots & \rho_j^{T-2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho_j^{T-1} & \rho_j^{T-2} & \rho_j^{T-3} & \dots & 1 \end{bmatrix}, \quad |\rho_i| < 1;$$

(3) For  $i, j = 1, 2, \dots, n$ :  $E \beta_i = \bar{\beta}$ .

$$E(\beta_i - \bar{\beta})(\beta_j - \bar{\beta})' = \begin{cases} \Delta & \text{if } i = j \\ 0 & \text{otherwise.} \end{cases}$$

$\Delta = \{\delta_{kk}\}$  ( $k, k' = 1, 2, \dots, K$ ) is positive definite;

(4)  $\beta_t$  is independent of  $u_j$  for  $i, j = 1, 2, \dots, n$ ;

(5) The  $x_{ikt}$  are exogenous variables distributed independently of  $\beta_t$  and  $u_t$ .<sup>2</sup> Furthermore,  $X_t$  is nonstochastic.

The implications of Assumption 1 are discussed by Swamy (1972). If we arrange the observations on each variable first by individual and then according to period, we may represent eq. (1) by

$$(2) \quad y = X\bar{\beta} + D_x \xi + u$$

where  $y \equiv (y_1', y_2', \dots, y_n')$ ,  $X \equiv [X_1', X_2', \dots, X_n']$ ,  $\bar{\beta} \equiv (\bar{\beta}_1, \bar{\beta}_2, \dots, \bar{\beta}_K)$ ,  $D_x \equiv \text{diag}[X_1, X_2, \dots, X_n]$ ,  $\xi \equiv [\xi_1', \xi_2', \dots, \xi_n']$ ,  $\beta_i = \bar{\beta} + \xi_i$ , and  $u \equiv (u_1', u_2', \dots, u_n')$ .

For given  $X$  the random vector  $y$  is distributed with mean  $X\bar{\beta}$  and variance-covariance (V-C) matrix of the form

$$(3) \quad \Sigma = \begin{bmatrix} X_1 \Delta X_1' + \sigma_{11} \Omega_{11} & \sigma_{12} \Omega_{12} & \dots & \sigma_{1n} \Omega_{1n} \\ \sigma_{21} \Omega_{21} & X_2 \Delta X_2' + \sigma_{22} \Omega_{22} & \dots & \sigma_{2n} \Omega_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} \Omega_{n1} & \sigma_{n2} \Omega_{n2} & \dots & X_n \Delta X_n' + \sigma_{nn} \Omega_{nn} \end{bmatrix}$$

<sup>1</sup> With an abuse of notation we use the same symbol to denote a random quantity and its value.

<sup>2</sup> This assumption is partly relaxed in Swamy (1972).

The objective is to estimate the parameter vector  $\theta = (\beta', \omega')'$  where  $\omega$  is a  $[n^2 + K^2 + n] \times 1$  vector containing all the elements of  $\{\sigma_{ij}\}$ ,  $\Delta$  and  $\rho_i$  ( $i = 1, 2, \dots, n$ ) arranged in any order.

Model (1) contains a sample space  $Y$  of elements  $y$ . The distribution of  $y$  over  $Y$  can be taken as known to belong to a continuously parameterized family of distributions with probability density function (pdf),  $p(y|X, \theta)$ , the parameter vector  $\theta$  ranging over a well-defined parameter space  $\Theta = \{\theta: -\infty < \beta_k < \infty, 0 < \delta_{kk} < \infty \text{ for } k = 1, 2, \dots, K; \delta_{kk}^2 < \delta_{kk}\delta_{k'k'}, \delta_{kk'} = \delta_{k'k} \text{ for } k \neq k' = 1, 2, \dots, K; \sigma_{ij}^2 < \sigma_{ii}\sigma_{jj}, \sigma_{ij} = \sigma_{ji} \text{ for } i \neq j = 1, 2, \dots, n; 0 < \sigma_{ii} < \infty, 0 \leq |\rho_i| < 1 \text{ for } i = 1, 2, \dots, n\}$ . We assume that the unknown true value of  $\theta$  belongs to  $\Theta$ .

## 2.2. Criteria of Estimation

Suppose that the seriousness of sampling errors,  $\hat{\beta} - \beta$ , is indicated by the loss matrix  $(\hat{\beta} - \beta)(\hat{\beta} - \beta)'$  and we wish to find an estimator  $\hat{\beta}$  for which

$$(4) \quad \mathbf{1}' E(\hat{\beta} - \beta)(\hat{\beta} - \beta) \mathbf{1}$$

is minimum for every  $\beta \in \Theta$  and every arbitrary vector  $\mathbf{1} \neq \mathbf{0}$ .

We assume that the loss matrix which expresses the demerit of the estimate  $\hat{\theta}$  of  $\theta$  is separable in its components  $\beta$  and  $\omega$ . We do not specify the loss function involving  $\omega$ . It is worth noting that in the problem of estimating  $\omega$  a quadratic loss function does not seem to be appropriate, see Ferguson (1967, p. 179). For each fixed  $\theta$ , the expected value of  $(\hat{\beta} - \beta)(\hat{\beta} - \beta)'$  relative to the distribution of  $y$  determined by  $\theta$  is called the risk matrix or the matrix of second order moments of  $\hat{\beta}$  around  $\beta$ .  $E(\hat{\beta}_k - \beta_k)^2$  is called the mean square error of  $\hat{\beta}_k$ .

A moment's reflection will reveal that it is not possible to find an estimator  $\hat{\beta}$  which minimizes (4) for every  $\beta \in \Theta$  and every  $\mathbf{1} \neq \mathbf{0}$ , see Silvey (1970, p. 24). For example, if we take  $\hat{\beta} = \mathbf{a}$  (a vector of constants) for all  $y$ , this estimator will have zero risk when  $\beta = \mathbf{a}$  and thus to have a better estimator in the sense of (4), an estimator  $\hat{\beta}$  must have zero risk for every  $\beta$ . This is obviously not possible. So we must modify our criterion of estimation.

As is well-known, if we restrict ourselves to a class of linear unbiased estimators of  $\beta$ , we can find an estimator which minimizes the risk in (4) for every  $\beta \in \Theta$  and every  $\mathbf{1} \neq \mathbf{0}$ . Such an estimator is the minimum variance linear unbiased (MVLU) estimator

$$(5) \quad \bar{b}(\omega) = (X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1}y.$$

In the practical situation in which  $\omega$  is unknown, an estimate  $\hat{\omega}$  of  $\omega$  developed by Swamy (1972) can be used in place of the known value used in (5). We can offer an asymptotic justification for this procedure.

It has been emphasized by many statisticians that there is an element of arbitrariness in the criterion of MVLU, particularly with regard to unbiasedness. Consequently, in what follows we modify the criterion of MVLU.

## 2.3. Stein-like Estimators

Following one approach of Zellner and Vandaele (1971), we consider the problem of estimating  $\beta$  when the loss function is quadratic. Specifically, let the

quadratic loss function be  $(\hat{\beta} - \beta)'Q(\hat{\beta} - \beta)$  where  $Q$  is a known positive definite matrix. Since the range of each  $\beta_k$  is bounded, the risk function  $E[(\hat{\beta} - \beta)'Q(\hat{\beta} - \beta)]$  is bounded, provided  $\hat{\beta}$  has finite  $V-C$  matrix. Zellner and Vandaale (1971) show that among all estimators of the form  $c\bar{\mathbf{b}}(\omega)$ , where  $c$  is a scalar lying between 0 and 1, the estimator

$$(6) \quad c^*\bar{\mathbf{b}}(\omega) \equiv \left[ 1 - \frac{\text{tr}(X'\Sigma^{-1}X)^{-1}Q}{\text{tr}(X'\Sigma^{-1}X)^{-1}Q + \beta'Q\beta} \right] \bar{\mathbf{b}}(\omega)$$

has the smallest risk. That is,  $E[c\bar{\mathbf{b}}(\omega) - \beta]'Q[c\bar{\mathbf{b}}(\omega) - \beta]$  takes the smallest value for every  $\beta \in \Theta$  when  $c = c^*$ . Also,

$$(7) \quad E[c^*\bar{\mathbf{b}}(\omega) - \beta]'Q[c^*\bar{\mathbf{b}}(\omega) - \beta] \leq E[\bar{\mathbf{b}}(\omega) - \beta]'Q[\bar{\mathbf{b}}(\omega) - \beta] \quad \forall \beta \in \Theta.$$

Since  $c^*\bar{\mathbf{b}}(\omega)$  involves parameters with unknown values, it cannot be computed. Therefore, as in Zellner and Vandaale (1971) we may approximate  $c^*\bar{\mathbf{b}}(\omega)$  by

$$(8) \quad \hat{c}^*\bar{\mathbf{b}}(\hat{\omega}) \equiv \left[ 1 - \frac{\text{tr}(X'\hat{\Sigma}^{-1}X)^{-1}Q}{\bar{\mathbf{b}}(\hat{\omega})'Q\bar{\mathbf{b}}(\hat{\omega})} \right] \bar{\mathbf{b}}(\hat{\omega})$$

where  $\hat{\Sigma}$  and  $\bar{\mathbf{b}}(\hat{\omega})$  are as shown in Swamy (1972).

The estimator  $\hat{c}^*\bar{\mathbf{b}}(\hat{\omega})$  is in the form of an estimator developed by Stein for the mean vector of a  $K$ -dimensional normal population, see Zellner and Vandaale (1971).

Following Mehta and Srinivasan (1971) we may approximate  $c^*$  by an exponential function with two adjustable parameters and write

$$(9) \quad \hat{f}(\gamma)\bar{\mathbf{b}}(\hat{\omega}) = [1 - \gamma_1 \exp\{-\gamma_2 \bar{\mathbf{b}}(\hat{\omega})'X'\hat{\Sigma}^{-1}X\bar{\mathbf{b}}(\hat{\omega})\}] \bar{\mathbf{b}}(\hat{\omega})$$

where  $0 < \gamma_1 < 1$  and  $\gamma_2 > 0$ .

Notice that the factor  $\hat{c}^*$  multiplying  $\bar{\mathbf{b}}(\hat{\omega})$  in (8) can take on negative values with positive probability. Baranchik's analysis of simpler situations (see Stein, 1966) indicates that the estimator in (8) can be improved upon by restricting  $\hat{c}^*$  to be nonnegative. The factor  $\hat{f}(\gamma)$  multiplying  $\bar{\mathbf{b}}(\hat{\omega})$  in (9) can be made positive by suitably choosing the values of  $\gamma_1$  and  $\gamma_2$ . Experience in simpler situations (Mehta and Srinivasan, 1971) has shown that by judicious choice of  $\gamma_1$  and  $\gamma_2$  one can make the risk associated with  $\hat{f}(\gamma)\bar{\mathbf{b}}(\hat{\omega})$  smaller than that associated with  $\bar{\mathbf{b}}(\hat{\omega})$  or with  $\hat{c}^*\bar{\mathbf{b}}(\hat{\omega})$  for a range of values of  $\beta$  around 0. Since the estimators in (8) and (9) provide only approximations to the optimal linear estimator  $c^*\bar{\mathbf{b}}(\omega)$ , neither of them is an estimator which has minimum average risk within the class of linear or nonlinear estimators of  $\beta$ , see Strawderman and Cohen (1971). Consequently, there are other ways of obtaining linear or nonlinear estimators which have smaller risks than  $\hat{c}^*\bar{\mathbf{b}}(\hat{\omega})$  and  $\hat{f}(\gamma)\bar{\mathbf{b}}(\hat{\omega})$  (see Section 4 below).

The estimator in (8) takes  $\bar{\mathbf{b}}(\hat{\omega})$  and pulls it towards a central value 0 or past 0 if  $\bar{\mathbf{b}}(\hat{\omega})'Q\bar{\mathbf{b}}(\hat{\omega}) < \text{tr}(X'\hat{\Sigma}^{-1}X)^{-1}Q$ .<sup>3</sup> Since all elements of  $\bar{\mathbf{b}}(\hat{\omega})$  are shrunk by the same factor towards 0, the extreme values experience most shift. The estimators in (8) and (9) may do very poorly in estimating those elements of  $\beta$  with unusually large

<sup>3</sup> If we knew *a priori* that the true values of the elements of  $\beta$  lay closely to a value other than zero, we could easily modify the formulae in (8) and (9) to shrink the estimated value of  $\beta$  towards that value, see Zellner and Vandaale (1971), and Mehta and Srinivasan (1971).

or small values. Unless the true values of all the elements of  $\beta$  lie closely in almost the same interval around 0, the estimators in (8) and (9) may not yield good estimates of all the elements of  $\beta$ . It may happen that for some values of  $\beta$  the total risk associated with (8) is smaller than that associated with  $\bar{b}(\hat{\omega})$  but the risk associated with an element of (8) is larger than that associated with the corresponding element of  $\bar{b}(\hat{\omega})$ . To put it differently, the estimator  $\hat{c}^*\bar{b}(\hat{\omega})$  may have good ensemble properties but not good component properties. This is also true of  $\hat{f}(\cdot)\bar{b}(\hat{\omega})$ .

To guard against this bad property of Stein-like estimators, Efron and Morris (1972) develop a "limited translation estimator" which is a compromise between Stein's estimator and the maximum likelihood estimator (MLE). The compromise consists of following the Stein rule as closely as possible subject to a fixed constraint on how far the estimator is allowed to deviate from MLE. This procedure is sensible if the probability that an ML estimator of  $\beta$  will be far removed from the true value of  $\beta$  is small. Indeed, this probability is large if  $X'\Sigma^{-1}X$  is close to singularity.

The average value of the squared distance from  $\bar{b}(\omega)$  to  $\beta$  is given by

$$(10) \quad E[\bar{b}(\omega) - \beta][\bar{b}(\omega) - \beta]' = \text{tr}(X'\Sigma^{-1}X)^{-1} = \sum_{i=1}^k \lambda_i^{-1}$$

where  $\lambda_i$  is a latent root of  $X'\Sigma^{-1}X$ . Consequently, if the set of independent variables is such that reasonable data collection results in an  $X'\Sigma^{-1}X$  with one or more latent roots close to 0, then the average distance from  $\bar{b}(\omega)$  to  $\beta$  will be large. In this case the Efron-Morris procedure of pulling an estimate of  $\beta$  towards  $\bar{b}(\omega)$  amounts to pulling an estimate away from  $\beta$ , which is not desirable. If the least squares estimates  $\bar{b}(\omega)$  lie far away from the true value of  $\beta$  as a result of high multicollinearity, then so will be the estimates given by  $\hat{f}(\cdot)\bar{b}(\hat{\omega})$  and  $\hat{c}^*\bar{b}(\hat{\omega})$ . Typically,  $X'X$  will not be close to a diagonal matrix in applications of economic relevance. In the next section we discuss procedures which are specifically designed to minimize the bad effects of significant departures of  $X'X$  from  $I$ . In order to guarantee good component properties we say that  $\hat{\beta}$  is "uniformly" better than  $\beta^*$  if

$$(11) \quad E(\hat{\beta} - \beta)(\hat{\beta} - \beta)' \leq E(\beta^* - \beta)(\beta^* - \beta)'$$

for every  $1 \neq 0$  and every  $\beta \in \Theta$ , with strict inequality for some  $\beta$ . In this way we avoid the specification of  $Q$ . An estimator,  $\hat{\beta}^*$ , is "inadmissible" if there exists another estimator of  $\beta$  which completely dominates  $\hat{\beta}^*$  in the sense of (11); otherwise it is "admissible". Notice that  $\hat{\beta}$  is uniformly better than  $\hat{\beta}^*$  in the sense of (11), if and only if,  $E(\hat{\beta}^* - \beta)(\hat{\beta}^* - \beta)'$  exceeds  $E(\hat{\beta} - \beta)(\hat{\beta} - \beta)'$  by a positive semi-definite matrix for every  $\theta \in \Theta$ .

### 3. SUGGESTED PROCEDURE OF ESTIMATION IN CASES OF PARTIAL PRIOR INFORMATION

#### 3.1. Ridge Regression

For the model in the present paper, let  $\hat{\beta}$  be constrained to be in a hypersphere of radius  $r$ . Let the estimation criterion be the minimum residual sum of squares  $(y - X\hat{\beta})\Sigma^{-1}(y - X\hat{\beta})$  subject to the condition  $\hat{\beta}'\hat{\beta} = r^2 < \infty$ . The value of  $\hat{\beta}$  that

minimizes the function

$$(12) \quad (y - X\bar{\beta})\Sigma^{-1}(y - X\bar{\beta}) + \mu(\bar{\beta}\bar{\beta} - r^2)$$

is

$$(13) \quad \bar{b}_\mu(\omega) = (X'\Sigma^{-1}X + \mu I)^{-1}X'\Sigma^{-1}y.$$

This is the ridge estimator developed by Hoerl and Kennard (1970a).

Unlike the Stein procedure, the above procedure takes into account the restrictions on the ranges of  $\bar{\beta}$ . The estimation procedure based on the matrix  $(X'\Sigma^{-1}X + \mu I)$  with  $\mu > 0$  rather than  $X'\Sigma^{-1}X$  can be used to circumvent many of the difficulties associated with the multicollinearity problem, and it can be used to obtain a point estimate of  $\bar{\beta}$ , which is on the average closer to  $\bar{\beta}$  than is  $\bar{b}(\omega)$ . The average value of the squared distance from  $\bar{b}_\mu(\omega)$  to  $\bar{\beta}$  is

$$(14) \quad E[\bar{b}_\mu(\omega) - \bar{\beta}][\bar{b}_\mu(\omega) - \bar{\beta}] = \text{tr} [I + \mu(X'\Sigma^{-1}X)^{-1}]^{-1}(X'\Sigma^{-1}X)^{-1} \\ \cdot [\mu(X'\Sigma^{-1}X)^{-1} + I]^{-1} + \mu^2\bar{\beta}'(X'\Sigma^{-1}X + \mu I)^{-2}\bar{\beta} \\ = \sum_{i=1}^k \lambda_i (\lambda_i + \mu)^2 + \mu^2\bar{\beta}'(X'\Sigma^{-1}X + \mu I)^{-2}\bar{\beta}.$$

This can be compared with (10). If a  $\lambda_i$  is close to zero, (14) will be substantially smaller than (10) depending on the value of  $\mu$ . That is, when  $X'\Sigma^{-1}X$  is ill-conditioned, the estimates of  $\bar{\beta}$  based on  $\bar{b}(\omega)$  (but not on  $\bar{b}_\mu(\omega)$ ) have a high probability of being far removed from  $\bar{\beta}$ . Hoerl and Kennard show that there exists a range of values of  $\mu$  for which the average distance from  $\bar{b}_\mu(\omega)$  to  $\bar{\beta}$  is smaller than that from  $\bar{b}(\omega)$  to  $\bar{\beta}$ .

The relationship of a ridge estimator to the Aitken estimator  $\bar{b}(\omega)$  is given by the alternative form

$$(15) \quad \bar{b}_\mu(\omega) = [I + \mu(X'\Sigma^{-1}X)^{-1}]^{-1}\bar{b}(\omega).$$

We may rewrite (13) as

$$(16) \quad \bar{b}_\mu(\omega) = (X'\Sigma^{-1}X + \mu I)^{-1}X'\Sigma^{-1}X\bar{b}(\omega).$$

The estimator  $\bar{b}_\mu(\omega)$  will be recognized as a "matrix weighted average" of the vectors  $\bar{b}(\omega)$  and  $0$ . Like  $c^*\bar{b}(\omega)$ , it also shrinks the estimated value of  $\bar{\beta}$  a fixed percentage away from  $\bar{b}(\omega)$  towards  $0$ . But the shrinkage factor is not the same for all the elements of  $\bar{b}(\omega)$ . Thus, the ridge regression technique, by utilizing the restriction on the range of  $\bar{\beta}$ , leads to an estimator which does not suffer from the limitations of  $c^*\bar{b}(\omega)$ . The estimator in (15) is insensitive to multicollinearity. On the other hand, when  $X'\Sigma^{-1}X = I$ , the matrix factor multiplying  $\bar{b}(\omega)$  in (16) reduces to a scalar times identity matrix. In this case, by appropriately defining  $\mu$  we can equate  $\bar{b}_\mu(\omega)$  to  $c^*\bar{b}(\omega)$ .

The second order moment matrix of  $\bar{b}_\mu(\omega)$  around  $\bar{\beta}$  is

$$(17) \quad E[\bar{b}_\mu(\omega) - \bar{\beta}][\bar{b}_\mu(\omega) - \bar{\beta}]' = [I + \mu(X'\Sigma^{-1}X)^{-1}]^{-1}(X'\Sigma^{-1}X)^{-1} \\ \cdot [I + \mu(X'\Sigma^{-1}X)^{-1}]^{-1} + \mu^2(X'\Sigma^{-1}X + \mu I)^{-1}\bar{\beta}\bar{\beta}'(X'\Sigma^{-1}X + \mu I)^{-1}.$$

The first term on the r.h.s. of (17) is the V-C matrix of  $\bar{b}_\mu(\omega)$  and the second term is the matrix of squares and cross products of the biases of the elements of  $\bar{b}_\mu(\omega)$ .

As is well-known,  $(X'\Sigma^{-1}X)^{-1}$  is the  $V-C$  matrix of  $\bar{\mathbf{b}}(\omega)$ . The matrix  $(X'\Sigma^{-1}X)^{-1} - [I + \mu(X'\Sigma^{-1}X)^{-1}]^{-1}(X'\Sigma^{-1}X)^{-1}[I + \mu(X'\Sigma^{-1}X)^{-1}]^{-1}$  is non-negative definite so that for some values of  $\mu$  and  $\bar{\beta}$  in a neighborhood of  $\mathbf{0}$  there is a possibility of  $E[\bar{\mathbf{b}}(\omega) - \bar{\beta}][\bar{\mathbf{b}}(\omega) - \bar{\beta}]' - E[\bar{\mathbf{b}}_{\mu}(\omega) - \bar{\beta}][\bar{\mathbf{b}}_{\mu}(\omega) - \bar{\beta}]'$  being positive semi-definite. However, the mean square error of an element of  $\bar{\mathbf{b}}_{\mu}(\omega)$  may not be substantially smaller than that of the corresponding element of  $\bar{\mathbf{b}}(\omega)$ , if the true value of  $\bar{\beta}$  is not sufficiently close to  $\mathbf{0}$ .

An approximate ridge regression estimator is

$$(18) \quad \bar{\mathbf{b}}_{\mu}(\hat{\omega}) = (X'\hat{\Sigma}^{-1}X + \mu I)^{-1}X'\hat{\Sigma}^{-1}\mathbf{y}.$$

In Hoerl and Kennard (1971b) some recommendations for choosing a  $\mu > 0$  are given.

### 3.2. Minimum Conditional Mean Square Error Estimator of $\bar{\beta}$

Recall that the second order moment matrix of a linear estimator  $A\mathbf{y} + \mathbf{a}$  around  $\bar{\beta}$  is

$$(19) \quad A\Sigma A' + [(AX - I)\bar{\beta} + \mathbf{a}][(AX - I)\bar{\beta} + \mathbf{a}]'.$$

The quantity in (19) cannot be minimized unless it is bounded, see Barnard (1963). Since the range of  $\theta$  is bounded, the elements of (19) are bounded. Let  $\bar{\beta}^*$  be a guessed value of  $\bar{\beta}$ . Using  $\bar{\beta}^*$  in place of  $\bar{\beta}$ , we obtain

$$(20) \quad A\Sigma A' + [(AX - I)\bar{\beta}^* + \mathbf{a}][(AX - I)\bar{\beta}^* + \mathbf{a}]'.$$

If (20) is chosen as a criterion of estimation, the optimum choice of  $\mathbf{a}$  is  $\mathbf{0}$  and that of  $A$  is (see Rao, 1971, p. 389)

$$(21) \quad A^* = \bar{\beta}^*\bar{\beta}^{*'}X'(X\bar{\beta}^*\bar{\beta}^{*'}X' + \Sigma)^{-1}.$$

Consequently, the optimal estimator of  $\bar{\beta}$ , given  $\bar{\beta}^*$ , is

$$(22) \quad \bar{\mathbf{b}}^*(\omega) = \bar{\beta}^*\bar{\beta}^{*'}X'(X\bar{\beta}^*\bar{\beta}^{*'}X' + \Sigma)^{-1}\mathbf{y}.$$

(Henceforth we shall refer to  $\bar{\mathbf{b}}^*(\omega)$  as the minimum conditional mean square error (MCMSE) estimator of  $\bar{\beta}$ . The result in (22) is given as an exercise in Theil (1971, p. 125, Problem 4.3.) Notice that the estimator  $\bar{\mathbf{b}}^*(\omega)$  exists even when the rank of  $X$  is less than  $K$ . In cases where the rank of  $X$  is  $K$ , we can write

$$(23) \quad (X\bar{\beta}^*\bar{\beta}^{*'}X' + \Sigma)^{-1} = \Sigma^{-1} - \Sigma^{-1}X(X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1} \\ + \Sigma^{-1}X(X'\Sigma^{-1}X)^{-1}[\bar{\beta}^*\bar{\beta}^{*'} + (X'\Sigma^{-1}X)^{-1}]^{-1} \\ \times (X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1},$$

see Rao (1965, p. 29, Problem 2.9).

Inserting this back into (22) gives

$$(24) \quad \bar{\mathbf{b}}^*(\omega) = \bar{\beta}^*\bar{\beta}^{*'}[\bar{\beta}^*\bar{\beta}^{*'} + (X'\Sigma^{-1}X)^{-1}]^{-1}\bar{\mathbf{b}}(\omega).$$

In the practical situation in which  $\Sigma$  is unknown, the estimator  $\bar{\mathbf{b}}^*(\omega)$  can be approximated by

$$(25) \quad \bar{\mathbf{b}}^*(\hat{\omega}) = \bar{\beta}^*\bar{\beta}^{*'}X'(X\bar{\beta}^*\bar{\beta}^{*'}X' + \hat{\Sigma})^{-1}\mathbf{y} \\ = \bar{\beta}^*\bar{\beta}^{*'}[\bar{\beta}^*\bar{\beta}^{*'} + (X'\hat{\Sigma}^{-1}X)^{-1}]^{-1}\bar{\mathbf{b}}(\hat{\omega}) \quad [\text{by (23)}]$$

where  $\hat{\Sigma}$  is as defined in (8).



In the Appendix to the paper, it will be shown that a sufficient condition for  $(X'\Sigma^{-1}X)^{-1} - E\{[\bar{\mathbf{b}}^*(\omega) - \beta][\bar{\mathbf{b}}^*(\omega) - \beta]'\beta^*\}$  to be positive definite is

$$(26) \quad \sup_k \bar{\beta}' \mathbf{p}_k \mathbf{p}_k' \bar{\beta} < 1$$

where  $\mathbf{p}_k$  is the  $k$ th column of  $P$ .  $P$  is a nonsingular matrix such that  $P(X'\Sigma^{-1}X)^{-1}P = I$ ,  $P\beta^*\beta^{*'}P = \lambda_1^2 \mathbf{i}_1 \mathbf{i}_1'$ , and  $\mathbf{i}_1$  is the first column of an identity matrix of order  $K$ .

It is clear from (A.4) in the Appendix that the conditional variance of an element of  $\bar{\mathbf{b}}^*(\omega)$ , given  $\bar{\beta}^*$ , is substantially smaller than the variance of the corresponding element of the Aitken estimator  $\bar{\mathbf{b}}(\omega)$  for every  $\theta$ . But, for some values of  $\theta$ , due to high magnitude of bias the conditional mean square error of an element of  $\bar{\mathbf{b}}^*(\omega)$ , given  $\bar{\beta}^*$ , exceeds the variance of the corresponding element of  $\bar{\mathbf{b}}(\omega)$ . Condition (26) indicates the values of  $\theta$  for which  $\bar{\mathbf{b}}^*(\omega)$  based on given  $\bar{\beta}^*$  is better than  $\bar{\mathbf{b}}(\omega)$ . Consequently, the approximate MCMSE estimator  $\bar{\mathbf{b}}^*(\hat{\omega})$  cannot completely dominate the approximate MVLU estimator  $\bar{\mathbf{b}}(\hat{\omega})$  in the sense of (11). When  $K = 1$ , condition (26) is satisfied if the square of the coefficient of variation of the MVLU estimator  $\bar{\mathbf{b}}(\omega)$  is greater than one. In the general case condition (26) is likely to be satisfied if  $X'\Sigma^{-1}X$  is close to singularity. Under these conditions, one can improve upon the MVLU estimator by relaxing the unbiasedness condition as in (20).

We now compare the moment matrices of  $\bar{\mathbf{b}}_\mu(\omega)$  and  $\bar{\mathbf{b}}^*(\omega)$ . It is seen from (A.4) and (A.6) in the Appendix that since the rank of  $\bar{\beta}^*\beta^{*'}$  is unity, the conditional variance of an element of  $\bar{\mathbf{b}}^*(\omega)$  is substantially smaller than the variance of the corresponding element of  $\bar{\mathbf{b}}_\mu(\omega)$ . However, for any reasonable values of  $\mu$  and  $\bar{\beta}^*$ , the magnitude of bias of an element of  $\bar{\mathbf{b}}^*(\omega)$  is likely to be larger than that of bias of the corresponding element of  $\bar{\mathbf{b}}_\mu(\omega)$ . For certain values of parameters,  $\bar{\mathbf{b}}^*(\hat{\omega})$  is better than  $\bar{\mathbf{b}}_\mu(\hat{\omega})$ .

Next, we note that, if a prior estimate of  $\bar{\beta}$  is not available, we may consider the following estimator:

$$(27) \quad \hat{\bar{\mathbf{b}}}(\hat{\omega}) = \bar{\mathbf{b}}_\mu(\hat{\omega})\bar{\mathbf{b}}'(\hat{\omega})X'[X\bar{\mathbf{b}}(\hat{\omega})\bar{\mathbf{b}}(\hat{\omega})X' + \hat{\Sigma}]^{-1}y.$$

When there is near-extreme multicollinearity, a precise estimation of  $\bar{\beta}$  is not possible, but a relatively precise estimation of  $X\bar{\beta}$  and  $\Sigma$  is possible, see Rao (1965, pp. 184-5) and Theil (1971, pp. 153-4). The estimator  $\hat{\bar{\mathbf{b}}}(\hat{\omega})$  is based on the precise estimates of  $\bar{\beta}$ ,  $X\bar{\beta}$  and  $\Sigma$ .

The estimator  $\bar{\mathbf{b}}^*(\omega)$  is based on a prior estimate of  $\bar{\beta}$ , while the estimator  $\bar{\mathbf{b}}_\mu(\omega)$  is based on a prior knowledge of the range of  $\bar{\beta}'\bar{\beta}$ . Since the rank of  $\bar{\beta}^*\beta^{*'}$  is unity, we cannot express  $\bar{\mathbf{b}}^*(\omega)$  in the form of a matrix weighted average of the vectors  $\bar{\mathbf{b}}(\omega)$  and  $\theta$ . However, when  $K = 1$ , by appropriately defining  $\mu$  we can equate  $\bar{\mathbf{b}}_\mu(\omega)$  to  $\bar{\mathbf{b}}^*(\omega)$ , see Theil (1971, p. 126, Problem 4.4).

In summary, we have found that none of the estimators  $\bar{\mathbf{b}}(\hat{\omega})$ ,  $\hat{c}^*\bar{\mathbf{b}}(\hat{\omega})$ ,  $\hat{f}(\cdot)\bar{\mathbf{b}}(\hat{\omega})$ ,  $\bar{\mathbf{b}}_\mu(\hat{\omega})$ , and  $\bar{\mathbf{b}}^*(\hat{\omega})$  is uniformly better than the other in the sense of (11). Consequently, it is not possible to choose among them unless we know "where in the parameter space to look" for the most efficient estimates. When we are faced with an extreme multicollinearity situation, we may use  $\bar{\mathbf{b}}^*(\hat{\omega})$  if a reliable prior estimate of  $\bar{\beta}$  is available and  $\hat{\bar{\mathbf{b}}}(\hat{\omega})$  otherwise.

#### 4. ESTIMATING PARAMETERS WITH THE FIRST TWO MOMENTS OF A PRIOR DISTRIBUTION

There are several situations in which extraneous information on some of the parameters of an equation is available. This information may arise from an analysis of past data and/or from theoretical and practical considerations; that is, from sources other than currently available sample. To incorporate such a prior information the following procedure was suggested by Durbin (1953) and developed further by Theil and Goldberger (1961) and Theil (1963).

##### 4.1. Mixed Estimation When $\bar{\beta}$ is Regarded as Fixed

Suppose that extraneous information of the following form is available.

$$(28) \quad \mathbf{r} = R\bar{\beta} + \mathbf{v} \quad \text{with} \quad E\mathbf{v} = \mathbf{0} \quad \text{and} \quad E\mathbf{v}\mathbf{v}' = \tau^2\psi,$$

where  $\mathbf{r}$  is a  $q \times 1$  vector of prior estimates of  $R\bar{\beta}$ .  $R$  is a  $q \times K$  matrix of known constants,  $\mathbf{v}$  is a  $q \times 1$  vector of errors in  $\mathbf{r}$  and  $q \leq K$ . We assume that  $\mathbf{v}$  is uncorrelated with  $\mathbf{u}$  and  $\xi$  in (2). We now combine equations (2) and (28) and apply the Aitken theorem to obtain the following estimator for  $\bar{\beta}$ .

$$(29) \quad \hat{\beta}_p(\omega) = \left( X'\Sigma^{-1}X + \frac{1}{\tau^2}R'\psi^{-1}R \right)^{-1} \left( X'\Sigma^{-1}\mathbf{y} + \frac{1}{\tau^2}R'\psi^{-1}\mathbf{r} \right).$$

The estimator  $\hat{\beta}_p(\omega)$  is the MVLU estimator of  $\bar{\beta}$  where linear now means linear in  $\mathbf{y}$  and  $\mathbf{r}$ . Here the distinction between  $\bar{\mathbf{b}}(\omega)$  as a MVLU estimator of  $\bar{\beta}$  and  $\hat{\beta}_p(\omega)$  as a MVLU estimator of the same  $\bar{\beta}$  is to be clearly understood. The linear function of  $\mathbf{y}$ , namely  $\bar{\mathbf{b}}(\omega)$ , is the MVLU estimator of  $\bar{\beta}$  in the sense that any other estimator of  $\bar{\beta}$  which is also linear in the vector  $\mathbf{y}$  and unbiased has a  $V$ - $C$  matrix which exceeds that of  $\bar{\mathbf{b}}(\omega)$  by a positive semidefinite matrix. On the other hand,  $\hat{\beta}_p(\omega)$  is the MVLU estimator of  $\bar{\beta}$  in the sense that any other estimator of  $\bar{\beta}$  which is linear in  $\mathbf{y}$  and  $\mathbf{r}$  and unbiased has a  $V$ - $C$  matrix which exceeds that of  $\hat{\beta}_p(\omega)$  by a positive semidefinite matrix. We shall refer to  $\hat{\beta}_p(\omega)$  as the "mixed regression" estimator. We again remind the reader that the criterion of MVLU is defective in its premises, in that the condition of unbiasedness sometimes leads to inadmissible estimates, see Ferguson (1967, pp. 135-6).

As  $\tau^2 \rightarrow 0$ , the estimator  $\hat{\beta}_p(\omega)$  approaches the restricted estimator of  $\bar{\beta}$  given by the normal equations (see Chipman, 1964, p. 1101)

$$(30) \quad \begin{bmatrix} X'\Sigma^{-1}X & R' \\ R & 0 \end{bmatrix} \begin{bmatrix} \bar{\beta} \\ \mu \end{bmatrix} = \begin{bmatrix} X'\Sigma^{-1}\mathbf{y} \\ \mathbf{r} \end{bmatrix}.$$

Eq. (30) is obtained by minimizing

$$(31) \quad \frac{1}{2}(\mathbf{y} - X\bar{\beta})\Sigma^{-1}(\mathbf{y} - X\bar{\beta}) - \mu'(\mathbf{r} - R\bar{\beta})$$

where  $\mu$  is a vector of Lagrangian multipliers. Theil and Goldberger (1961), solve eq. (30) under the assumption that the ranks of  $X$  and  $R$  are  $K$  and  $q$  respectively, while Rao and Mitra (1971, p. 147) solve the same equation without any restrictions on the ranks of  $X$  and  $R$ .

Chipman (1964, pp. 1101-2) points out an important special case of (29). If  $\psi$  is known, eq. (28) can be written as

$$(32) \quad \begin{aligned} \psi^{-1/2} \mathbf{r} &= \psi^{-1/2} R \hat{\boldsymbol{\beta}} + \psi^{-1/2} \mathbf{v} \\ \mathbf{r}^* &= R^* \hat{\boldsymbol{\beta}} + \mathbf{v}^* \end{aligned}$$

When the rank of  $X$  is less than  $K$ , and when  $X^* = X \Sigma^{-1/2}$  and  $R^*$  are "complementary",<sup>4</sup>  $X^* = (X \Sigma^{-1} X + (1/\tau^2) R \psi^{-1} R)^{-1} X^*$  is a generalized inverse of  $X^*$ , independently of  $1/\tau^2$ , as long as  $0 < 1/\tau^2 < \infty$ , because for all such  $1/\tau^2$ ,  $1/\tau R^*$  has the same row space as  $R^*$ . Similarly,  $R^* = (\tau^2 X \Sigma^{-1} X + R \psi^{-1} R)^{-1} R^*$  is a generalized inverse of  $R^*$ , independently of  $1/\tau^2$ , as long as  $0 < \tau^2 < \infty$ , since for all such  $\tau^2$ ,  $\tau X^*$  has the same row space as  $X^*$ . Therefore the estimator  $\hat{\boldsymbol{\beta}}_p(\omega)$  is functionally independent of  $1/\tau^2$  as long as  $0 < 1/\tau^2 < \infty$  and  $R^*$  is complementary to  $X^*$ . In this case the estimator  $\hat{\boldsymbol{\beta}}_p(\omega)$  can be computed even when  $\tau^2$  is unknown.

To consider another case, let  $q = K$  and  $R = I$ . Then  $\hat{\boldsymbol{\beta}}_p(\omega)$  becomes

$$(33) \quad \hat{\boldsymbol{\beta}}_p(\omega) = \left( X \Sigma^{-1} X + \frac{1}{\tau^2} \psi^{-1} \right)^{-1} \left( X \Sigma^{-1} \mathbf{y} + \frac{1}{\tau^2} \psi^{-1} \mathbf{r} \right).$$

It is easily seen that  $\hat{\boldsymbol{\beta}}_p(\omega)$  in (33) is a "matrix weighted average" of  $\bar{\mathbf{b}}(\omega)$  and  $\mathbf{r}$ , with weights inversely proportional to their respective  $F$ - $C$  matrices. Hence, an estimate of  $\bar{\boldsymbol{\beta}}$  is pulled towards  $\mathbf{r}$  away from  $\bar{\mathbf{b}}(\omega)$ . The estimator in (33) covers  $\bar{\mathbf{b}}_p(\omega)$  in (13) as a special case. When  $\mathbf{r} = \mathbf{0}$  and  $\tau^2 \psi = (1/\mu)I$ ,  $\bar{\mathbf{b}}_p(\omega)$  is the same as (33).

Analysis of simpler situations has shown that the estimator

$$(34) \quad \hat{\boldsymbol{\beta}}_p(\hat{\omega}) = \left( X \hat{\Sigma}^{-1} X + \frac{1}{\tau^2} \psi^{-1} \right)^{-1} \left( X \hat{\Sigma}^{-1} \mathbf{y} + \frac{1}{\tau^2} \psi^{-1} \mathbf{r} \right),$$

with known  $\tau^2 \psi$ , completely dominates  $\bar{\mathbf{b}}(\hat{\omega})$  in the sense of (11), provided  $E\mathbf{v} = \mathbf{0}$  and  $\xi$  and  $\mathbf{u}$  are normal, see Swamy and Mehta (1969), and Mehta and Swamy (1972b). In cases where  $E(\mathbf{v} - \boldsymbol{\eta})(\mathbf{v} - \boldsymbol{\eta})' = \tau^2 \psi$ ,  $\boldsymbol{\eta}$  is unknown,  $\tau^2 \psi$  is known, and  $\xi$  and  $\mathbf{u}$  are normal,  $\hat{\boldsymbol{\beta}}_p(\hat{\omega})$  is better than  $\bar{\mathbf{b}}(\hat{\omega})$  if only the coefficient of variation of each element of  $\mathbf{v}$  is sufficiently large in magnitude, see Swamy and Mehta (1972). Thus, if we misspecify the prior moments, there is no guarantee that each diagonal element of the second order moment matrix of  $\hat{\boldsymbol{\beta}}_p(\hat{\omega})$  around  $\bar{\boldsymbol{\beta}}$  will be less than or equal to the corresponding diagonal element of the second order moment matrix of  $\bar{\mathbf{b}}(\hat{\omega})$  around  $\bar{\boldsymbol{\beta}}$ .

The compatibility test statistic developed by Theil (1963) can be utilized to test whether prior information is in conflict with sample information. Mehta and Swamy (1972a) have derived the exact finite sample distribution of Theil's compatibility test statistic. They have also considered the consequences for estimation, in terms of mean square error, of making preliminary tests. The efficiency of preliminary testing procedures has been examined by comparison of the risk functions of preliminary test estimators with that of pure regression estimator,  $\bar{\mathbf{b}}(\omega)$ , which is an Aitken estimator when no prior information is used. The preliminary test estimator dominated the pure regression estimator over certain regions of the parameter space.

<sup>4</sup> The matrices  $X^*$  and  $R^*$  are complementary if (1)  $\text{rank}(X^*) + \text{rank}(R^*) = K$ , (2)  $X^*$  and  $R^*$  have the same number of columns, and (3) the row spaces of  $X^*$  and  $R^*$  have only the origin in common.

Returning again to the case where  $E\mathbf{v} = \mathbf{0}$  and  $E\mathbf{v}\mathbf{v}' = \tau^2\psi$ , it can be seen that the matrix

$$\begin{aligned} & \left( X'\Sigma^{-1}X + \frac{1}{\tau^2}\psi^{-1} \right)^{-1} - \beta^*\beta^{*'}[\beta^*\beta^{*'} + (X'\Sigma^{-1}X)^{-1}]^{-1}(X'\Sigma^{-1}X)^{-1} \\ & \times [\bar{\beta}^*\bar{\beta}^{*'} + (X'\Sigma^{-1}X)^{-1}]^{-1}\bar{\beta}^*\bar{\beta}^{*'} - (X'\Sigma^{-1}X)^{-1}[\bar{\beta}^*\bar{\beta}^{*'} \\ & + (X'\Sigma^{-1}X)^{-1}]^{-1}\bar{\beta}\bar{\beta}'[\bar{\beta}^*\bar{\beta}^{*'} + (X'\Sigma^{-1}X)^{-1}]^{-1}(X'\Sigma^{-1}X)^{-1} \end{aligned}$$

is positive definite only for certain values of  $\theta$ ,  $\beta^*$  and  $\tau^2\psi$ . Consequently, the estimator  $\bar{b}^*(\hat{\omega})$  in (25) will not be uniformly better than  $\hat{\beta}_p(\hat{\omega})$  in (34) even when the first two moments of  $\mathbf{r}$  are exactly known.

A particular case which can be solved exactly, and for which there is a complete and simpler treatment is the following. Let  $K = 1$ , and  $\beta^{*2} = |r^2 - \tau^2\psi|$ . Notice that  $E r^2 = \beta^2 + \tau^2\psi$ . We can use standard analytical and numerical methods (Mehta and Swamy, 1972a) to evaluate the unconditional mean square error of  $\bar{b}^*(\omega)$  with respect to the distributions of  $\beta^{*2}$  and  $\mathbf{y}$ . If the square of the coefficient of variation of  $r$ ,  $(\tau^2\psi/\beta^2)$ , is greater than one and the square of the coefficient of variation of the MVLU estimator  $\bar{b}(\omega)$  is greater than or equal to one, then  $\bar{b}^*(\omega)$  is better than  $\hat{\beta}_p(\omega)$ .

Formulae (25) and (34) provide two different ways of combining prior information with sample information. Neither one of them is better than the other regardless of the true values of parameters. It should be emphasized that the estimator  $\bar{b}^*(\hat{\omega})$  should not be used unless  $\beta^*\beta^{*'}$  is a reliable estimate of  $\bar{\beta}\bar{\beta}'$ . If the prior point estimates of the elements of  $\bar{\beta}$  are not reliable, then it is better to express the uncertainties associated with these estimates in the form of a distribution with mean  $\bar{\beta}$  and  $V-C$  matrix  $\tau^2\psi$  and use the estimator  $\hat{\beta}_p(\hat{\omega})$ . That the prior information be unbiased is a severe restriction on the nature of such information, see Zellner (1970, p. 189). This restriction will be eliminated in the next subsection.

#### 4.2. Bayesian Estimation When $\bar{\beta}$ is Regarded As a Random Variable

We now make the following "wide-sense" assumption.

*Assumption 2:* A probability distribution on a class of measurable sets in  $\Theta$  exists. The variable  $\bar{\beta}$  is judged *a priori* to be distributed independently of  $\omega$  whose distribution is a point distribution with the whole mass of the distribution concentrated at one point. Furthermore,  $E\bar{\beta} = \mathbf{r}$  and  $E(\bar{\beta} - \mathbf{r})(\bar{\beta} - \mathbf{r})' = \tau^2\psi$  which is positive definite.

Even if a purely pragmatic attitude is adopted it does seem to be true that for at least some inference problems, an approach which assumes the existence of a prior distribution of  $\theta$  is more appropriate than one which does not. However, it is very restrictive to assume that the distribution of  $\omega$  is a point distribution. If this assumption is relaxed, the analysis gets very complicated, see Lindley and Smith (1972).

Assuming that  $\theta$  is a random variable, Zellner and Vandaele (1971) discuss the Bayesian interpretations (attributable to Lindley and others) of the Stein-like estimator  $c^*\bar{b}(\omega)$ . When  $X'X = I$ ,  $Q = I$ ,  $\Sigma = \sigma^2I$ , and the prior distribution of  $\bar{\beta}$  has mean  $\theta$  and scalar  $V-C$  matrix, one can generate a Bayes estimator of the form

$c^*\bar{b}(\omega)$ . Notice that when  $r$  is regarded as a fixed parameter,  $\hat{\beta}_p(\omega)$  is still a linear function of  $y$  but becomes a biased estimator of  $\beta$ . It is interesting to note that if  $\beta$  is considered to be a random variable with mean equal to fixed  $r$  and fixed  $V-C$  matrix  $\tau^2\psi$ , then  $\hat{\beta}_p(\omega)$  in (33) is the "best linear" predictor of  $\beta$  in the sense that any other predictor of  $\beta$  which is also linear in the vector  $y$  has an averaged second order moment matrix around  $\beta$  which exceeds that of  $\hat{\beta}_p(\omega)$  by a positive semidefinite matrix. In other words, if  $r$  and  $\tau^2\psi$  are the mean and  $V-C$  matrix of  $\beta$ , then  $\hat{\beta}_p(\omega)$  completely dominates every other linear in  $y$  estimator (predictor) of  $\beta$  in the sense of (11).<sup>5</sup> Proof of this important result is given in Chipman (1964, p. 1105) and Rao (1965, p. 192). If  $r \neq 0$ ,  $\psi \neq I$  and Assumption 2 is true, the formulae  $c^*\bar{b}(\omega)$  and  $\bar{b}_p(\omega)$  are inappropriate. When  $\beta$  is random, the procedure outlined in subsection 3.2 is also inappropriate because, under Assumption 2, (19) is not the second order moment matrix of  $Ay + a$  around  $\beta$ , see Chipman (1964, p. 1104). Notice that the estimators  $\bar{b}(\omega)$ ,  $c^*\bar{b}(\omega)$ ,  $\bar{b}_p(\omega)$  and  $\bar{b}^*(\omega)$  for given  $\beta^*$ , are all linear functions of  $y$ . Hence, it follows from the Chipman-Rao theorem that they are inferior to the best linear estimator  $\hat{\beta}_p(\omega)$  if Assumption 2 is true. Thus, the biased estimators generated through the Chipman-Rao procedure are better than those generated through the procedure outlined in subsection 3.2.

We called  $\hat{\beta}_p(\omega)$  the best linear estimator of  $\beta$ . The qualification linear can be dropped if the prior distribution of  $\beta$ , given  $r$  and  $\tau^2\psi$ , is normal and the conditional distribution of  $y$ , given  $X$ ,  $\Sigma$ , and  $\beta$ , is also normal. This is because, under these normality assumptions, the estimator  $\hat{\beta}_p(\omega)$  is the mean of the conditional posterior distribution of  $\beta$ , given  $\Sigma$ ,  $r$ ,  $\tau^2\psi$  and the data, see Zellner (1971, p. 76), and Zellner and Vandaele (1971). The posterior mean  $\hat{\beta}_p(\omega)$  with known  $\Sigma$ ,  $\tau^2\psi$  and  $r$  is admissible with respect to a quadratic loss function, see Zellner (1971, p. 24). Thus, admissible estimates can be found if the prior distribution of  $\theta$  is completely known, see Ferguson (1967).

Even though the result in (33) is intuitively appealing, it has certain weaknesses. In (13) and (33) different posterior means have been obtained by combining two different priors with the same likelihood of parameters. These priors were therefore influential in deciding the posterior means in small samples. It is worth noting that if the Aitken estimate  $\bar{b}(\omega)$  and the prior mean  $r$  are very different, then the estimate (33) is a long way from  $\bar{b}(\omega)$ . In this case it may happen that either the model specification is at fault or the prior information is incompatible with sample information, see Box and Jenkins (1970, p. 251). Efron and Morris (1971) also point out that the estimator  $\hat{\beta}_p(\omega)$  must give bad estimates when  $r$  is far from  $\beta$ . Let  $N_K(r, \tau^2\psi)$  represent the true prior distribution of  $\beta$ .<sup>6</sup> Suppose that this distribution is actually a mixture of various other distributions, one of which is  $N_K(r_1, \tau_1^2\psi_1)$  such that  $\tau^2\psi - \tau_1^2\psi_1$  is positive definite. For any fixed value of  $\tau_1^2\psi_1$ , the expected squared error risk of an element of  $\hat{\beta}_p(\omega)$  with respect to the prior distribution  $N_K(r_1, \tau_1^2\psi_1)$  can be made arbitrarily large by moving  $r_1$  arbitrarily far from  $r$ . That is, the estimator  $\hat{\beta}_p(\omega)$  does well on the population,  $N_K(r, \tau^2\psi)$ , as a whole, but may perform very poorly on a particular subpopulation,  $N_K(r_1, \tau_1^2\psi_1)$ . The estimator  $(X'\Sigma^{-1}X + (1/\tau_1^2)\psi_1^{-1})^{-1}(X'\Sigma^{-1}y + (1/\tau_1^2)\psi_1^{-1}r_1)$

<sup>5</sup> The requirement that an estimator of  $\beta$  be linear arises from the absence, in our "distribution-free" formulation, of the assumption about the form of the prior distribution of  $\beta$ .

<sup>6</sup>  $N_K(r, \tau^2\psi)$  represents  $K$ -dimensional normal with mean  $r$  and  $V-C$  matrix  $\tau^2\psi$ .

does well on the subpopulation  $N_k(\mathbf{r}_1, \tau_1^2 \psi_1)$ . If we knew that a particular  $\beta$  belonged to the subpopulation  $N_k(\mathbf{r}_1, \tau_1^2 \psi_1)$ , then we could use the estimator  $(X' \Sigma^{-1} X + (1/\tau_1^2) \psi_1^{-1})^{-1} (X' \Sigma^{-1} y + (1/\tau_1^2) \psi_1^{-1} \mathbf{r}_1)$  rather than  $\hat{\beta}_p(\omega)$ . Information on subpopulation distributions can be obtained by assessing  $\mathbf{r}$  and  $\tau^2 \psi$  as precisely as possible. Now the relevant question is: How can we assess a prior distribution in practice?

Notice that the probability distribution on a class of measurable sets in  $\Theta$  is viewed merely as a reflection of the belief of the statistician about where the true value of  $\theta$  lies prior to an observation being made. Conditions under which such a distribution exists are given in Ferguson (1967, Section 1.4). It has been shown by Savage and others that personal probabilities assessed in accordance with certain plausible behavioral postulates of "coherence" must conform mathematically to a probability measure, see Lindley (1971). Winkler (1967a,b; 1971) discusses the practical problem of the assessment of personal probabilities. An operational way of assessing a probability is through the study of relevant gambles. Methods such as scoring rules and bets are useful in leading individuals to make careful probability assessments.

It should be emphasized, however, that in many economic situations there remains the practical difficulty of assessing a prior distribution to reflect one's degree of belief. If the parameter space contains a finite number of points, then by sufficient introspection one can arrive at the prior odds at which one would just accept a bet on this parameter value rather than that, and so eventually find the prior distribution appropriate for a particular problem. If  $\Theta$  is continuous, as it usually is, it is not clear whether any reasonable consideration of the way in which inferences cohere leads to the existence of the prior distribution, see Lindley (1971, pp. 7-8). The difficulty of choosing a prior distribution is highlighted, when the parameter space is infinite-dimensional as in Sims (1971). Efron and Morris (1971, p. 808) argue that in the realistic situations there is seldom any one prior distribution that is "true" in an absolute sense. There are only more or less relevant priors. If a distribution with mean  $\mathbf{r}$  and  $V$ - $C$  matrix  $\tau^2 \psi$  is at all in doubt, it would be well to modify the estimator  $\hat{\beta}_p(\omega)$ .

In large samples the situation improves. With a reasonably informative experiment, the values  $\mathbf{r}$  and  $\tau^2 \psi$  adequate for describing rather imprecise knowledge can be changed quite considerably without affecting the final result all that much. This is the consequence of the fact that, under general conditions, sample information dominates prior information in fairly large samples. In fact, Lindley (1971, p. 62) has shown that if the pdf  $p(y|X, \theta)$  satisfies certain regularity conditions (see Silvey, 1961 and Perlman, 1972), the method of maximum likelihood is shown to be a reasonably "coherent" technique in large samples. We, therefore, turn to a study of this topic.

## 5. MAXIMUM LIKELIHOOD METHOD

In this section we assume the following:

*Assumption 3:* Given  $X$ ,  $\hat{\beta}$ , and  $\Sigma$ ,  $y$  is normally distributed with mean  $X\hat{\beta}$  and  $V$ - $C$  matrix  $\Sigma$ , i.e.,  $y \sim N_{nT}(X\hat{\beta}, \Sigma)$ .

For simplicity, we let  $\sigma_{ij} = 0$  if  $i \neq j$  and  $\rho_i = 0$  for every  $i$ . Now  $\theta = (\bar{\beta}', \omega')$  where  $\omega$  is a  $(n + K^2) \times 1$  vector.  $\omega$  denotes the vector presentation of the  $\sigma_{ii}$ 's and all elements of  $\Delta$  in which  $\sigma_{11}, \dots, \sigma_{nn}$  appear in order first, then the elements of the first column of  $\Delta$ , the elements of the second column and so on.

The pdf of  $y$ , given  $X$ , is

$$(35) \quad p(y|X, \theta) = (2\pi)^{-nT/2} \prod_{i=1}^n \left\{ \sigma_{ii}^{-1(T-K)} |X_i'X_i|^{-1/2} |\Delta + \sigma_{ii}(X_i'X_i)^{-1}|^{-1/2} \right\} \\ \cdot \exp \left\{ -\frac{1}{2} \sum_{i=1}^n \left[ \frac{(T-K)s_{ii}}{\sigma_{ii}} + (\mathbf{b}_i - \bar{\beta})' \right. \right. \\ \left. \left. \cdot [\Delta + \sigma_{ii}(X_i'X_i)^{-1}]^{-1} (\mathbf{b}_i - \bar{\beta}) \right] \right\}$$

where

$$s_{ii} = y_i' M_i y_i / (T - K), \quad M_i = I - X_i(X_i'X_i)^{-1}X_i'$$

and

$$\mathbf{b}_i = (X_i'X_i)^{-1}X_i'y_i,$$

see Swamy (1971, pp. 111-12).

Now, given the data  $y, X$ ,  $p(y|X, \theta)$  in (35) may be regarded as a function of  $\theta$ . When so regarded, it is called the likelihood function of  $\theta$  for given  $y$  and  $X$ . The likelihood function is defined up to a multiplicative constant. The likelihood expresses the relative plausibilities of different parameter values after we have observed the data  $y$  and  $X$ , see Barnard (1967). Methods of eliminating nuisance parameters from the likelihood function so that inferences can be made about the parameters of interest are considered by Kalbfleisch and Sprott (1970). In this regard "marginal" and "conditional" likelihoods are introduced. These can be computed if only the likelihood function factors into two parts, one of which contains a parameter of interest, say  $\beta_k$ , only and the other being uninformative about  $\beta_k$  in the absence of knowledge of other parameters. It is clear from (35) that the likelihood function has the form (apart from irrelevant constants)

$$(36) \quad l(\theta|y, X) \propto \left[ \prod_{i=1}^n \sigma_{ii}^{-1(T-K)} \right] \exp \left\{ -\frac{1}{2} \sum_{i=1}^n \frac{(T-K)s_{ii}}{\sigma_{ii}} \right\} \\ \cdot \left[ \prod_{i=1}^n |\Delta + \sigma_{ii}(X_i'X_i)^{-1}|^{-1/2} \right] \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (\mathbf{b}_i - \bar{\beta})' [\Delta \right. \\ \left. + \sigma_{ii}(X_i'X_i)^{-1}]^{-1} (\mathbf{b}_i - \bar{\beta}) \right\}.$$

Each of the first  $n$  factors on the right hand side of (36) contains one of the  $\sigma_{ii}$  only. It contains no available information concerning  $\bar{\beta}$  and  $\Delta$  in the absence of

knowledge of the  $\sigma_{ii}$ . Unfortunately, the last factor contains available information about every element of  $\theta$ , see Kalbfleisch and Sprott (1970, p. 200). However, as  $T \rightarrow \infty$  since  $(X_i'X_i)^{-1} \rightarrow 0$ , the last factor gives less and less information about the  $\sigma_{ii}$ 's.  $\beta$  and  $\Delta$  are the parameters of our interest and we cannot derive their marginal likelihoods from (36). It is meaningless to integrate  $l(\theta|y, X)$  in an attempt to obtain the marginal likelihoods of the elements of  $\beta$ , see Box and Tiao (1973, p. 73). However, a close study of the likelihood function is always desirable. In certain instances, the data will contain no information regarding certain parameters. It is important to study the likelihood function's properties to determine when this is the case, see, for example, Box and Jenkins (1970, pp. 225-6), Silvey (1970, pp. 81-2), Swamy and Mehta (1971), and Swamy and Rao (1971). A general method for obtaining a reasonable estimate of  $\theta$  in most situations is the well-known maximum likelihood method, see Rao (1962). In this section we try to verify the conditions which ensure the consistency and asymptotic normality of an ML estimator  $\hat{\theta}$  of  $\theta$ . First, we indicate a method of obtaining  $\hat{\theta}$ .

An ML estimate  $\hat{\theta}$  is any element of  $\Theta$  such that  $p(y|X, \hat{\theta}) = \sup_{\theta \in \Theta} p(y|X, \theta)$ .  $\hat{\theta}$  belongs to the set which is most plausible after we have observed  $y$  and  $X$ . At this point it should be appreciated that the ML method always estimates the entire underlying distribution from given data. Successful estimation of the entire underlying distribution is the maximum of objectives attainable by any statistical method. Since  $\Theta$  is an open set, it may happen that no ML estimate of  $\theta$  exists. However, a neighborhood ML estimate of  $\theta$ , which is defined by Kiefer and Wolfowitz (1956, p. 892), exists in some cases where an ML estimate does not. Usually, ML estimates emerge as a solution of the likelihood equations  $\partial \log l(\theta|y, X)/\partial \theta = 0$  shown in Swamy (1971, p. 112). These equations are nonlinear in the unknowns and have to be solved numerically. A convenient method of solving the likelihood equations is the method of scoring described in Rao (1965, p. 302), see also Silvey (1970, 70-1). This method requires an explicit derivation of information matrix which is given by (see Swamy, 1971, p. 114)

$$(36) \quad I(\theta) = \begin{bmatrix} \frac{E\hat{c}^2 \log l}{\partial \beta \partial \beta'} & \frac{E\hat{c}^2 \log l}{\partial \beta \partial \omega'} \\ & \frac{E\hat{c}^2 \log l}{\partial \omega \partial \omega'} \end{bmatrix}$$

where

$$\begin{aligned} \frac{E\hat{c}^2 \log l}{\partial \beta \partial \beta'} &= \sum_{i=1}^n [\Delta + \sigma_{ii}(X_i'X_i)^{-1}]^{-1}, & \frac{E\hat{c}^2 \log l}{\partial \beta \partial \omega'} &= 0, \\ \frac{E\hat{c}^2 \log l}{\partial \omega \partial \omega'} &= \begin{bmatrix} \left\{ \frac{E\hat{c}^2 \log l}{\partial \sigma_{ii} \partial \sigma_{jj}} \right\} & \left\{ \frac{E\hat{c}^2 \log l}{\partial \sigma_{ii} \partial \Delta_r} \right\} \\ & \left\{ \frac{E\hat{c}^2 \log l}{\partial \Delta_r \partial \Delta_r} \right\} \end{bmatrix}. \end{aligned}$$



$\Delta$ , denotes the vector presentation of all elements of  $\Delta$  in which the elements of the first row appear in order first, then the elements of the second row and so on:

$$\begin{aligned} -\frac{E\hat{c}^2 \log 1}{\partial\sigma_{ii}\partial\sigma_{ii}} &= \frac{1}{2} \frac{(T-K)}{\sigma_{ii}^2} + \frac{1}{2} \text{tr} [\Delta + \sigma_{ii}(X_i'X_i)^{-1}]^{-1} (X_i'X_i)^{-1} \\ &\quad \cdot [\Delta + \sigma_{ii}(X_i'X_i)^{-1}]^{-1} (X_i'X_i)^{-1} \quad (i = 1, 2, \dots, n) \\ -\frac{E\hat{c}^2 \log 1}{\partial\sigma_{ii}\partial\sigma_{jj}} &= 0 \quad \text{if } i \neq j, \\ -\frac{E\hat{c}^2 \log 1}{\partial\sigma_{ii}\partial\Delta} &= \frac{1}{2} [\Delta + \sigma_{ii}(X_i'X_i)^{-1}]^{-1} (X_i'X_i)^{-1} [\Delta + \sigma_{ii}(X_i'X_i)^{-1}]^{-1} \\ &\quad (i = 1, 2, \dots, n) \\ -\frac{E\hat{c}^2 \log 1}{\partial\Delta_r\partial\Delta_r} &= \frac{1}{2} \sum_{i=1}^n [\Delta + \sigma_{ii}(X_i'X_i)^{-1}]^{-1} \otimes [\Delta + \sigma_{ii}(X_i'X_i)^{-1}]^{-1}. \end{aligned}$$

$\otimes$  denotes the Kronecker product, see Tracy and Dwyer (1969, pp. 1580, 88-89).

### 5.1. Consistency of An ML Estimator of $\theta$

The pdf  $p(\mathbf{y}|X, \theta)$  in (35) depends on an unknown parameter vector  $\theta$  belonging to a metric space  $\Theta$  which is a subset of  $[K + n + \frac{1}{2}K(K + 1)]$ -dimensional Euclidean space. In (35) there is a family of possible distributions given by different values of  $\theta$  in  $\Theta$  and we do not know which one is appropriate. Let  $\theta_0$  be the unknown true value of  $\theta$ . We shall denote by  $E_0 \log p(\mathbf{y}|X, \theta)$  and  $\text{var}_0 \log p(\mathbf{y}|X, \theta)$  the mean and variance respectively of the random variable  $\log p(\mathbf{y}|X, \theta)$  on the sample space  $Y$  (of elements  $\mathbf{y}$ ) with respect to the distribution of  $\mathbf{y}$  determined by  $\theta_0$ . Let  $N_0$  be an open neighborhood of  $\theta_0$ . To prove that  $\hat{\theta}$  is weakly consistent we have to show that  $[\log p(\mathbf{y}|X, \theta_0) - \sup_{\theta \in \Theta - N_0} \log p(\mathbf{y}|X, \theta)] > 0$  in probability according to  $p(\mathbf{y}|X, \theta_0)$  see Silvey (1961, pp. 445-6). This means that the value of  $\theta$  which maximizes  $l(\theta|\mathbf{y}, X)$  belongs to  $N_0$  in probability when  $\theta_0$  obtains. If for every  $n, T$  and  $\theta \neq \theta_0$ , we have  $E_0 \log p(\mathbf{y}|X, \theta_0) > E_0 \log p(\mathbf{y}|X, \theta)$ , and  $E_0 \{\log p(\mathbf{y}|X, \theta_0) - \log p(\mathbf{y}|X, \theta)\}$  is large relative to  $[\text{var}_0 \{\log p(\mathbf{y}|X, \theta_0) - \log p(\mathbf{y}|X, \theta)\}]^{1/2}$ , then it follows from Chebychev's inequality that the method of maximum likelihood will discriminate well between  $\theta_0$  and other  $\theta$ . By putting certain regularity conditions on  $l(\theta|\mathbf{y}, X)$  we can guarantee that the method will discriminate well between  $\theta_0$  and, simultaneously, all other parameter values outside an open neighborhood of  $\theta_0$ , for large enough  $n$  and  $T$ . This is the basis of consistency proofs given by Silvey and others.

The likelihood function in (36) contains terms of different orders, each containing a particular subvector of  $\theta$ . Consequently, we proceed as follows: First, we assume that  $\theta_0 \in \Theta$ . Second, we rewrite (35) as

$$(37) \quad p(\mathbf{y}|X, \theta) = \left[ \prod_{i=1}^n g(s_{ii}|\sigma_{ii}) \right] f(\mathbf{b}|X, \theta)$$

where

$$g(s_{ii}|\sigma_{ii})z\sigma_{ii}^{-(T-K)} \exp \left\{ -\frac{1}{2} \frac{(T-K)s_{ii}}{\sigma_{ii}} \right\}$$

and

$$f(\mathbf{b}|X, \boldsymbol{\theta})z \left[ \prod_{i=1}^n [\Delta + \sigma_{ii}(X_i'X_i)^{-1}]^{-1/2} \right] \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (\mathbf{b}_i - \boldsymbol{\beta})' \cdot [\Delta + \sigma_{ii}(X_i'X_i)^{-1}]^{-1} (\mathbf{b}_i - \boldsymbol{\beta}) \right\}.$$

By Jensen's inequality (Silvey, 1970, p. 75) we have

$$(38) \quad E_0 \left[ \frac{1}{nT} \sum_{i=1}^n \log g(s_{ii}|\sigma_{ii0}) + \frac{1}{n} \log f(\mathbf{b}|X, \boldsymbol{\theta}_0) \right] \geq E_0 \left[ \frac{1}{nT} \sum_{i=1}^n \log g(s_{ii}|\sigma_{ii}) + \frac{1}{n} \log f(\mathbf{b}|X, \boldsymbol{\theta}) \right]$$

where  $\sigma_{ii0}$  is the true value of  $\sigma_{ii}$ . The inequality in (38) is strict unless  $\boldsymbol{\theta} = \boldsymbol{\theta}_0$ , because, in view of Assumption 1,  $\boldsymbol{\theta}$  is identified and the distributions corresponding to  $\boldsymbol{\theta}_0$  and  $\boldsymbol{\theta}$  are different.

There is a connection between "local" identifiability of a vector-valued parameter  $\boldsymbol{\theta}$  and positive definiteness of the information matrix  $I(\boldsymbol{\theta})$ , see Rothenberg (1971) and Silvey (1970, pp. 81-2).

*Assumption 4:* The vectors  $\mathbf{x}_{it} = (x_{i1t}, x_{i2t}, \dots, x_{iKt})'$  are all contained in a compact subset of  $K$ -dimensional Euclidean space such that for each  $i = 1, 2, \dots, n$  the matrix  $T^{-1}X_i'X_i$  converges to a finite positive definite matrix as  $T \rightarrow \infty$ .

Let  $D = \text{diag}[nI_K, TI_n, nI_{K^2}]$ . Now consider  $D^{-1/2}I(\boldsymbol{\theta}_0)D^{-1/2}$  where  $I(\boldsymbol{\theta}_0)$  is obtained from (36) by replacing  $\boldsymbol{\theta}$  by  $\boldsymbol{\theta}_0$ . The positive definiteness of  $\lim_{T \rightarrow \infty, n \rightarrow \infty} D^{-1/2}I(\boldsymbol{\theta}_0)D^{-1/2}$  which is necessary for the local identifiability of  $\boldsymbol{\theta}_0$  follows from Assumption 4. Following the same argument as in Silvey (1970, pp. 81-2) we can show that for any  $\boldsymbol{\theta} \neq \boldsymbol{\theta}_0$

$$(39) \quad \lim_{\substack{T \rightarrow \infty \\ n \rightarrow \infty}} E_0 \left[ \frac{1}{nT} \sum_{i=1}^n \log g(s_{ii}|\sigma_{ii0}) + \frac{1}{n} \log f(\mathbf{b}|X, \boldsymbol{\theta}_0) \right] > \lim_{\substack{T \rightarrow \infty \\ n \rightarrow \infty}} E_0 \left[ \frac{1}{nT} \sum_{i=1}^n \log g(s_{ii}|\sigma_{ii}) + \frac{1}{n} \log f(\mathbf{b}|X, \boldsymbol{\theta}) \right].$$

It is easy to show that for every  $\boldsymbol{\theta} \in \Theta$

$$(40) \quad E_0 \left[ \frac{1}{nT} \sum_{i=1}^n \log g(s_{ii}|\sigma_{ii}) + \frac{1}{n} \log f(\mathbf{b}|X, \boldsymbol{\theta}) \right] = O(1),$$

$$(41) \quad \text{var}_0 \left[ \frac{1}{nT} \sum_{i=1}^n \log g(s_{ii}|\sigma_{ii}) \right] = O\left(\frac{1}{nT}\right),$$

and

$$(42) \quad \text{var}_0 \left[ \frac{1}{n} \log f(\mathbf{b}|X, \boldsymbol{\theta}) \right] = O(n^{-1}).$$

Let  $\Theta_0$  be a compact subset of  $\Theta$ .

Assumption 5:  $\theta_0 \in \Theta_s$ .

In various practical situations it is often possible to rule out sufficiently extreme values of  $\theta$  on theoretical grounds and form  $\Theta_s$ , so that  $\theta_0 \in \Theta_s$ . In cases where the maximum likelihood procedure outlined in the previous subsection leads to implausible estimates like negative estimates for the diagonal elements of  $\Delta$ , Assumption 5 may not hold. In these cases we should examine Assumption 1 more closely. Under certain additional conditions we can replace Assumption 5 by a wider condition, see Perlman (1972).

The function  $p(y|X, \theta)$  is a pdf on the sample space  $Y$ , given  $X$ , for each  $\theta$  in  $\Theta_s$ , and the function  $l(\theta|y, X)$  is continuous on the metric space  $\Theta_s$  for each  $y$ , given  $X$ . Since  $\Theta_s - N_0$  is compact, we can cover it by a finite number, say  $h$ , of open spheres of radius  $r_{N_0}$ , having centers  $\theta_1, \dots, \theta_h$ , say. Let  $\log p(y|X, \theta_m, r_{N_0})$  be the supremum of  $\log p(y|X, \theta_j)$  with respect to  $\theta_j$  when  $\|\theta_m - \theta_j\| < r_{N_0}$ . For any  $\theta_m \in \Theta_s$  we have,  $\lim E_0 \log p(y|X, \theta_m, r_{N_0}) < \infty$ , as  $r_{N_0} \rightarrow 0$  because  $p(y|X, \theta)$  is uniformly bounded in  $y, \theta$  and  $E_0 \log p(y|X, \theta_0) < \infty$ . We can show that

$$(43) \quad E_0 \left[ \frac{1}{nT} \sum_{i=1}^n \log g(s_{ii}|\sigma_{ii}, r_{N_0}) + \frac{1}{n} \log f(\mathbf{b}|X, \theta_m, r_{N_0}) \right] \\ < E_0 \left[ \frac{1}{nT} \sum_{i=1}^n \log g(s_{ii}|\sigma_{ii}, r_{N_0}) + \frac{1}{n} \log f(\mathbf{b}|X, \theta_0) \right] \quad (m = 1, 2, \dots, h).$$

The results in (38)–(43) are adequate to establish the consistency of an ML estimate of  $\theta$ , see Swamy and Rao (1971), and Silvey (1961).

## 5.2. Asymptotic Normality

The standard method of establishing the asymptotic normality of an ML estimator  $\hat{\theta}$  of  $\theta$  utilizes the following results:

(a) Taylor's theorem in the expansion of  $\partial \log l(\hat{\theta}|y, X)/\partial \theta_0$ ;

(b) a central limit theorem applied to  $D^{-1/2}(\partial \log l(\theta_0|y, X)/\partial \theta_0)$ ;

(c) a law of large numbers applied to  $D^{-1/2}(\partial^2 \log l(\theta_0|y, X)/\partial \theta_0 \partial \theta_0)D^{-1/2}$ .

Under Assumptions 1, 3, 4 and 5 we have enough regularity conditions to establish the above results, see Silvey (1971, pp. 77–8) and Swamy and Rao (1971). Consequently,  $D^{-1/2}(\hat{\theta} - \theta_0)$  is asymptotically normal with mean  $\mathbf{0}$  and V-C matrix  $[\lim_{T \rightarrow \infty} \lim_{n \rightarrow \infty} D^{-1/2} I(\theta_0) D^{-1/2}]^{-1}$ .

The argument just presented, combined with the fact that the prior distribution of  $\theta$  does not depend on  $n$  and  $T$ , shows that in large samples, when Assumptions 1–5 are satisfied, the posterior distribution of  $\theta$  is approximately normal with mean  $\hat{\theta}$  and V-C matrix  $[-(\partial^2 \log l(\theta|y, X)/\partial \theta \partial \theta)]^{-1}$  evaluated at  $\hat{\theta}$ ; see Lindley (1971, p. 62) and Zellner (1971, pp. 32–3). This result is true even when the prior distribution of  $\omega$  is not a point distribution, provided the above conditions are satisfied.

## 6. SUMMARY AND CONCLUSIONS

In this paper we considered six different estimators of the mean of a random coefficient vector. These are (1) the MVLU estimator  $\bar{\mathbf{b}}(\omega)$ , (2) the Stein-like estimator  $c^* \bar{\mathbf{b}}(\omega)$ , (3) the ridge regression estimator  $\bar{\mathbf{b}}_\mu(\omega)$ , (4) the MCMSE

estimator  $\bar{b}^*(\omega)$ , (5) the mixed regression estimator  $\hat{\beta}_p(\omega)$ , and (6) an ML estimator  $\hat{\beta}$  of  $\beta$ . We also found feasible approximations to these estimators. None of the estimators  $\bar{b}(\omega)$ ,  $c^*\bar{b}(\omega)$ ,  $\bar{b}_\mu(\omega)$  and  $\bar{b}^*(\omega)$  is uniformly better than the other. Each of these estimators has its own weaknesses. In cases where *a priori* unbiased estimator  $r$  of  $\beta$  is available and its  $V$ - $C$  matrix  $\tau^2\psi$  is known, the estimator  $\hat{\beta}_p(\omega)$  is uniformly better than the estimator  $\bar{b}(\omega)$ . Under these conditions, the estimator  $\hat{\beta}_p(\omega)$  is also better than  $\bar{b}^*(\omega)$  if  $\bar{\beta}^*\bar{\beta}^*$  is not a reliable estimate of  $\beta\beta'$ . The estimators  $\bar{b}_\mu(\omega)$ ,  $\bar{b}^*(\omega)$  and  $\hat{\beta}_p(\omega)$  are insensitive to extreme multicollinearity. The estimator  $\hat{\beta}_p(\omega)$  covers the estimators  $c^*\bar{b}(\omega)$  and  $\bar{b}_\mu(\omega)$  as special cases.

When  $\beta$  is regarded as a random variable, the formula  $\bar{b}^*(\omega)$  is inappropriate and the estimator  $\hat{\beta}_p(\omega)$  covers the estimators  $\bar{b}_\mu(\omega)$  and  $c^*\bar{b}(\omega)$  as special cases. The prior information utilized in obtaining the estimator  $\hat{\beta}_p(\omega)$  is likely to provide a better numerical approximation to the practical situation than those utilized in obtaining the estimators  $c^*\bar{b}(\omega)$  and  $\bar{b}_\mu(\omega)$ . The estimator  $\hat{\beta}_p(\omega)$  is uniformly better than the estimators  $\bar{b}(\omega)$ ,  $c^*\bar{b}(\omega)$ ,  $\bar{b}_\mu(\omega)$  and  $\bar{b}^*(\omega)$  if  $\beta$  is distributed with mean  $r$  and  $V$ - $C$  matrix  $\tau^2\psi$ . Furthermore,  $\hat{\beta}_p(\omega)$  has all the desirable properties of a posterior mean corresponding to a normal prior and normal likelihood. In small samples one cannot find a uniformly better estimator of  $\beta$  unless the prior distribution of  $\beta$  is proper and known.

Under certain regularity conditions, the maximum likelihood estimate  $\hat{\beta}$  is at least as good as any other estimator of  $\beta$  in large samples.

Federal Reserve System, Washington D.C.

#### APPENDIX

Here we provide the proof of (26). The conditional second order moment matrix of  $\bar{b}^*(\omega)$  in (24) around  $\bar{\beta}$ , given  $\bar{\beta}^*$ , is

$$(A.1) \quad \bar{\beta}^*\bar{\beta}^*[\bar{\beta}^*\bar{\beta}^* + (X'\Sigma^{-1}X)^{-1}]^{-1}(X'\Sigma^{-1}X)^{-1}[\bar{\beta}^*\bar{\beta}^* + (X'\Sigma^{-1}X)^{-1}]^{-1}\bar{\beta}^*\bar{\beta}^* \\ + (X'\Sigma^{-1}X)^{-1}[\bar{\beta}^*\bar{\beta}^* + (X'\Sigma^{-1}X)^{-1}]^{-1}\bar{\beta}\bar{\beta}'[\bar{\beta}^*\bar{\beta}^* \\ + (X'\Sigma^{-1}X)^{-1}]^{-1}(X'\Sigma^{-1}X)^{-1}.$$

The first term in (A.1) is the conditional  $V$ - $C$  matrix of  $\bar{b}^*(\omega)$  and the second term is the matrix of squares and cross-products of the biases of the elements of  $\bar{b}^*(\omega)$  for given  $\bar{\beta}^*$ . Subtracting (A.1) from the  $V$ - $C$  matrix of  $\bar{b}(\omega)$  gives

$$(A.2) \quad (X'\Sigma^{-1}X)^{-1} - \bar{\beta}^*\bar{\beta}^*[\bar{\beta}^*\bar{\beta}^* + (X'\Sigma^{-1}X)^{-1}]^{-1}(X'\Sigma^{-1}X)^{-1} \\ - [\bar{\beta}^*\bar{\beta}^* + (X'\Sigma^{-1}X)^{-1}]^{-1}\bar{\beta}^*\bar{\beta}^* - (X'\Sigma^{-1}X)^{-1}[\bar{\beta}^*\bar{\beta}^* + (X'\Sigma^{-1}X)^{-1}]^{-1} \\ \cdot \bar{\beta}\bar{\beta}'[\bar{\beta}^*\bar{\beta}^* + (X'\Sigma^{-1}X)^{-1}]^{-1}(X'\Sigma^{-1}X)^{-1}.$$

Let  $P$  be a nonsingular matrix such that  $P'(X'\Sigma^{-1}X)^{-1}P = I$ , and  $P\bar{\beta}^*\bar{\beta}^*P = \lambda_1^*i_1i_1'$  where  $i_1$  is the first column of an identity matrix of order  $K$ . We pre and post multiply (A.2) by  $P^{-1}P'$  and  $PP^{-1}$  respectively to obtain

$$(A.3) \quad P^{-1}P^{-1} - P^{-1}\{\lambda_1^*i_1i_1'(\lambda_1^*i_1i_1' + I)^{-2}\lambda_1^*i_1i_1'\}P^{-1} \\ - P^{-1}(\lambda_1^*i_1i_1' + I)^{-1}\lambda_1^*0_10_1'(\lambda_1^*i_1i_1' + I)^{-1}P^{-1}$$

where  $\mathbf{0}_1$  is the characteristic vector corresponding to the nonzero root  $\lambda_1$  of  $P\bar{\beta}\beta'P$ . Using an identity in Swamy (1971, p. 25, Lemma 2.2.2) we have

$$(A.4) \quad P^{-1}P^{-1} - P^{-1}\mathbf{i}_1\mathbf{i}_1'P^{-1} \frac{\lambda_1^{*2}}{(1 + \lambda_1^*)^2} = P^{-1} \left( I - \frac{\lambda_1^*}{1 + \lambda_1^*} \mathbf{i}_1\mathbf{i}_1' \right) \lambda_1 \mathbf{0}_1 \mathbf{0}_1' \\ \cdot \left( I - \frac{\lambda_1^*}{1 + \lambda_1^*} \mathbf{i}_1\mathbf{i}_1' \right) P^{-1}.$$

Consequently, given  $\bar{\beta}^*$ ,

$$(A.5) \quad E[\bar{\mathbf{b}}(\omega) - \bar{\beta}][\bar{\mathbf{b}}(\omega) - \bar{\beta}]' = E[\bar{\mathbf{b}}^*(\omega) - \bar{\beta}][\bar{\mathbf{b}}^*(\omega) - \bar{\beta}]' \\ = P^{-1} \left\{ I - \mathbf{i}_1\mathbf{i}_1' \left[ \frac{\lambda_1^{*2}}{(1 + \lambda_1^*)^2} + \frac{\lambda_1^{*2} \lambda_1 \mathbf{0}_{11}^2}{(1 + \lambda_1^*)^2} \right] - \lambda_1 \mathbf{0}_1 \mathbf{0}_1' + \frac{\lambda_1^* \lambda_1 \mathbf{0}_{11}}{(1 + \lambda_1^*)} \right. \\ \left. \cdot (\mathbf{i}_1 \mathbf{0}_1' + \mathbf{0}_1 \mathbf{i}_1') \right\} P^{-1}$$

where  $\mathbf{0}_{11}$  is the first element of  $\mathbf{0}_1$ .

Let the matrix within the curl brackets be  $B$ . The matrix in (A.5) is positive definite if  $B$  is positive definite. Since  $B$  is symmetric,  $B$  is positive definite if all its diagonal elements are positive. The first diagonal element of  $B$  is positive if  $\bar{\beta}'\mathbf{p}_1\mathbf{p}_1'\bar{\beta} < 1 + 2\lambda_1^*$  where  $\mathbf{p}_1$  is the first column of  $P$ . Every other diagonal element of  $B$  is positive if  $\bar{\beta}'\mathbf{p}_k\mathbf{p}_k'\bar{\beta} < 1$   $k = 2, \dots, K$ .

Using  $P$  we may rewrite (17) as

$$(A.6) \quad P^{-1} \{ P'P(P'P + \mu I)^{-2} P'P \} P^{-1} + P^{-1} \{ \mu^2 (P'P + \mu I)^{-1} P\bar{\beta}\beta'P \\ \cdot (P'P + \mu I)^{-1} \} P^{-1}.$$

#### REFERENCES

- Barnard, G. (1963), "The Logic of Least Squares," *Journal of the Royal Statistical Society, Series B*, 25, pp. 124-7.
- (1967), "The Use of the Likelihood Function in Statistical Practice," pp. 27-40, in *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, 1, Berkeley and Los Angeles: University of California Press.
- Box, G. E. P. and G. M. Jenkins (1970), *Time Series Analysis, Forecasting, and Control*, San Francisco: Holden-Day.
- Box, G. E. P. and G. C. Tiao (1973), *Bayesian Inference in Statistical Analysis*, Massachusetts: Addison-Wesley Publishing Company.
- Chipman, J. S. (1964), "On Least Squares with Insufficient Observations," *Journal of the American Statistical Association*, 59, pp. 1078-111.
- Durbin, J. (1953), "A Note on Regression When There is Extraneous Information About one of the Coefficients," *Journal of the American Statistical Association*, 48, pp. 799-808.
- Efron, B. and C. Morris (1971), "Limiting the Risk of Bayes and Empirical Bayes Estimators. Part I: The Bayes Case," *Journal of the American Statistical Association*, 66, pp. 807-15.
- (1972), "Limiting Risk of Bayes and Empirical Bayes Estimators—Part II: The Empirical Bayes Case," *Journal of the American Statistical Association*, 67, pp. 130-9.
- Feige, E. L. and P. A. V. B. Swamy (1972), "A Random Coefficient Model of the Demand for Liquid Assets," Workshop Paper # 7211, Social Systems Research Institute, The University of Wisconsin, Madison. Accepted for publication in *Journal of Money, Credit, and Banking*.
- Ferguson, T. S. (1967), *Mathematical Statistics*, New York: Academic Press.
- Hoerl, A. E. and R. W. Kennard (1970a), "Ridge Regression: Biased Estimation for Non-Orthogonal Problems," *Technometrics*, 12, pp. 55-67.
- (1970b), "Ridge Regression: Applications to Non-Orthogonal Problems," *Technometrics*, 12, pp. 69-82.

- Kalbfleisch, J. G. and D. A. Sprott (1970). "Application of Likelihood Methods to Models Involving Large Numbers of Parameters." *Journal of the Royal Statistical Society, Series B*, 32, pp. 175-94.
- Kiefer, J. and J. Wolfowitz (1956). "Consistency of the Maximum Likelihood Estimator in the Presence of Infinitely Many Nuisance Parameters." *Annals of Mathematical Statistics*, 27, pp. 887-906.
- Lindley, D. V. (1971). *Bayesian Statistics, A Review*. SIAM, Philadelphia.
- Lindley, D. V. and A. F. M. Smith (1972). "Bayes Estimates for the Linear Model." *Journal of the Royal Statistical Society, Series B*, 34, pp. 1-41.
- Mehta, J. S. and R. Srinivasan (1971). "Estimation of the Mean by Shrinkage to a Point." *Journal of the American Statistical Association*, 66, pp. 86-90.
- Mehta, J. S. and P. A. V. B. Swamy (1972a). "The Exact Finite Sample Distribution of Theil's Compatibility Statistic and its Application." Report 7204. Division of Economic Research, Department of Economics, Ohio State University, Columbus. Accepted for publication in *Journal of the American Statistical Association*.
- . (1972b). "Efficient Method of Estimating the Level of a Stationary First-Order Autoregressive Process." *Communications in Statistics* (forthcoming).
- Perlman, M. D. (1972). "On the Strong Consistency of Approximate Maximum Likelihood Estimators." pp. 263-81 in *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability*, I. Berkeley and Los Angeles: California University Press.
- Rao, C. R. (1962). "Apparent Anomalies and Irregularities in Maximum Likelihood Estimation (with discussion)." *Sankhyā, Series A*, 24, pp. 73-102.
- . (1965). *Linear Statistical Inference and its Applications*. New York: John Wiley & Sons.
- . (1971). "Unified Theory of Linear Estimation." *Sankhyā, Series A*, 33, pp. 371-94.
- Rao, C. R. and S. K. Mitra (1971). *Generalized Inverse of Matrices and its Applications*. New York: John Wiley & Sons.
- Rothenberg, T. J. (1971). "Identification in Parametric Models." *Econometrica*, 39, pp. 577-92.
- Silvey, S. D. (1961). "A Note on Maximum-Likelihood in the Case of Dependent Random Variables." *Journal of the Royal Statistical Society, Series B*, 23, pp. 444-52.
- . (1970). *Statistical Inference*. Baltimore: Penguin.
- Sims, C. A. (1971). "Distributed Lag Estimation When the Parameter Space is Explicitly Infinite-Dimensional." *Annals of Mathematical Statistics*, 42, pp. 1622-36.
- Stein, C. (1966). "An Approach to the Recovery of Inter-Block Information in Balanced Incomplete Block Designs." pp. 351-66. in F. N. David (editor), *Research Papers in Statistics*. Festschrift for J. Neyman. New York: John Wiley & Sons.
- Strawderman, W. E. and A. Cohen (1971). "Admissibility of Estimators of the Mean Vector of a Multivariate Normal Distribution with Quadratic Loss." *Annals of Mathematical Statistics*, 42, pp. 720-96.
- Swamy, P. A. V. B. (1971). *Statistical Inference in Random Coefficient Regression Models*. New York: Springer-Verlag.
- . (1972). "Linear Models with Random Coefficients." to appear in Paul Zarembka (editor), *Frontiers of Econometrics*. New York: Academic Press, forthcoming.
- Swamy, P. A. V. B. and J. S. Mehta (1969). "On Theil's Mixed Regression Estimator." *Journal of the American Statistical Association*, 64, pp. 273-6.
- . (1972). "Robustness of Theil's Mixed Regression Estimators." Report 7218. Division of Economic Research, Department of Economics, Ohio State University, Columbus.
- . (1971). "Bayesian Analysis of Error Components Regression Models." Report 7203. Division of Economic Research, Department of Economics, Ohio State University, Columbus. to appear in *Journal of the American Statistical Association*.
- Swamy, P. A. V. B. and J. N. K. Rao (1971). "On Consistency and Asymptotic Normality of Maximum Likelihood Estimators of Parameters in a Distributed Lag Model with Auto-correlated Errors." Technical Report #16. University of Manitoba, Winnipeg.
- Theil, H. (1963). "On the Use of Incomplete Prior Information in Regression Analysis." *Journal of the American Statistical Association*, 58, pp. 401-14.
- . (1971). *Principles of Econometrics*. New York: John Wiley & Sons.
- Theil, H. and A. S. Goldberger (1961). "On Pure and Mixed Statistical Estimation in Economics." *International Economic Review*, 2, pp. 65-78.
- Tracy, D. S. and P. S. Dwyer (1969). "Multivariate Maxima and Minima with Matrix Derivatives." *Journal of the American Statistical Association*, 64, pp. 1576-94.
- Winkler, R. L. (1967a). "The Assessment of Prior Distributions in Bayesian Analysis." *Journal of the American Statistical Association*, 62, pp. 776-800.
- . (1967b). "The Quantification of Judgment: Some Methodological Suggestions." *Journal of the American Statistical Association*, 62, pp. 1109-20.
- . (1971). "Probabilistic Prediction: Some Experimental Results." *Journal of the American Statistical Association*, 66, pp. 675-85.

- Zellner, A. (1970), "The Bayesian Approach and Alternatives to Econometrics - I," pp. 178-93 in M. D. Intriligator (editor), *Frontiers of Quantitative Economics*. Amsterdam: North-Holland Publishing Company.
- (1971), *An Introduction to Bayesian Inference in Econometrics*. New York: John Wiley & Sons.
- Zellner, A. and W. Vandaale (1971), "Bayes-Stein Estimators for K-Means, Regression and Simultaneous Equation Models," presented at the Third NBER-NSF Symposium on Bayesian Inference in Econometrics.