



The World's Largest Open Access Agricultural & Applied Economics Digital Library

This document is discoverable and free to researchers across the globe due to the work of AgEcon Search.

Help ensure our sustainability.

Give to AgEcon Search

AgEcon Search

<http://ageconsearch.umn.edu>

aesearch@umn.edu

*Papers downloaded from **AgEcon Search** may be used for non-commercial purposes and personal study only. No other use, including posting to another Internet site, is permitted without permission from the copyright owner (not AgEcon Search), or as allowed under the provisions of Fair Use, U.S. Copyright Act, Title 17 U.S.C.*

No endorsement of AgEcon Search or its fundraising activities by the author(s) of the following work or their employer(s) is intended or implied.

AMSTER

AE 7/85

FACULTY OF
ACTUARIAL SCIENCE
&
ECONOMETRICS

GIANNINI FOUNDATION OF
AGRICULTURAL ECONOMICS
LIBRARY

OCT 10 1985

A & E REPORT

REPORT AE 7/85

ON MULTIVARIATE RIDGE REGRESSION

Y. Haitovsky (Hebrew University of Jerusalem)



University of Amsterdam

Title: On multivariate ridge regression

Author: Y. Haitovsky (Hebrew University of Jerusalem)

Address: University of Amsterdam
Faculty of Actuarial Science and Econometrics
Jodenbreestraat 23
1011 NH Amsterdam

Date: April 1985

Series and Number: Report AE 7/85

Pages: 28

Price: No charge

JEL Subject Classification: 211

Keywords: multivariate regression; ridge regression, Bayes estimates; empirical Bayes; least squares; maximum likelihood; linear hierarchical model; Stein's estimator; mean of a multivariate normal distribution; Exchangeability; election night forecasting.

Abstract

A multivariate linear regression model with q responses as a linear function of p independent variables $\mathbf{Y} = \mathbf{XB} + \mathbf{E}$ is considered with a $p \times q$ parameter matrix \mathbf{B} . The least squares (or Maximum Likelihood for multivariate normal \mathbf{E}) estimator of \mathbf{B} is deficient in that it takes no account of the "across regression" correlations, on the one hand, and ignores the famous Stein effect, on the other hand. A remedy was offered by Brown and Zidek (1980) in the form of a multivariate ridge estimator. A richer class of estimators is obtained here by casting the model in a linear hierarchical framework, obtaining the Brown and Zidek multivariate ridge estimators., Efron and Morris' estimators of several normal mean vectors and Fearn's Bayesian estimators of growth curves as special cases. The unknown covariance cases result in an identifiability problem which is treated in a Bayesian fashion using conjugate priors. The method is then applied to forecasting the final election results from partial returns obtained at election night.

On Multivariate Ridge Regression

1. INTRODUCTION:

Brown and Zidek (1980) consider a multivariate problem with q responses and n observations, \mathbf{Y} , assumed to satisfy the standard multivariate linear regression model

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E} \quad (1)$$

with \mathbf{X} an $(n \times p)$ full column rank matrix whose elements

are treated as fixed known constants

and \mathbf{B} is a $(p \times q)$ matrix of unknown coefficients to be estimated. With $\mathbf{E} = (\epsilon_1, \dots, \epsilon_q)$, the usual assumptions on the error are

$$E(\epsilon_j) = 0, \text{cov}(\epsilon_j, \epsilon_l) = \gamma_{jl}I_n \quad j, l = 1, \dots, q \quad (2)$$

or, in short $E(\mathbf{E}) = 0$, $\text{cov}(\mathbf{E}) = \Gamma \otimes I_n$, where \otimes denotes the usual Kronecker product of matrices. The least squares estimator of \mathbf{B} is given by

$$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \quad (3)$$

which is also the maximum likelihood estimator when normality of the error

distribution is assumed. Equivalently, writing $\hat{\mathbf{B}} = (\hat{\beta}_1, \dots, \hat{\beta}_q)$,

$\mathbf{Y} = (y_1, \dots, y_q)$ so that the vectors y_j ($n \times 1$) and $\hat{\beta}_j$ ($p \times 1$) pertain to the j -th of the q responses, (3) can be rewritten as

$$\hat{\beta}_j = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y_j \quad (j = 1, \dots, q) \quad (3a)$$

Brown and Zidek cite Sclove's (1971) argument that the least squares estimator (3) is deficient in that it takes no account of $\Gamma = (\gamma_{ij})$, the between regressions covariance matrix. In order to remedy this shortcoming and also to

take advantage of possible improvements in the estimation of the β_i 's even when $\Gamma = I_\gamma$ by using the ridge regression estimator for that case, Brown and Zidek suggest the multivariate ridge regression estimator (MRRE) of the form

$$\hat{\beta}^*(K) = (I_{\gamma_q} \otimes X^T X + K \otimes I_p)^{-1} (I_{\gamma_q} \otimes X^T) y \quad (4)$$

where $K(q \times q) > 0$ is the ridge matrix, $y = \text{vec}\{Y\}$ is the $(nq \times 1)$ vector

$y^T = (y_{11}^T, \dots, y_{\gamma q}^T)$ and $\hat{\beta}^*(K)$ is a $(pq \times 1)$ vector of estimators of $\beta = \text{vec}\{\beta\} = (\beta_1^T, \dots, \beta_{\gamma q}^T)^T$ where each β_j is a $(p \times 1)$ column vector.¹

The least squares estimators(3) have also a stacked matrix representation

$$\hat{\beta} = (I_{\gamma_q} \otimes X^T X)^{-1} (I_{\gamma_q} \otimes X^T) y. \quad (3b)$$

Clearly $\hat{\beta}^*(0) = \hat{\beta}$.

Estimator (4) is recognized by Brown and Zidek as a Bayes estimator, which can be directly obtained from Lindley and Smith's (1972) treatment of the linear model. However, a more careful application of the Lindley-Smith approach will yield a richer class of estimators, which is the aim of this paper.

It will be shown that the relationship between the resulting estimator and the Brown-Zidek MRRE is similar to that between the Lindley-Smith Exchangeability Within Regression estimator and the Ridge Regression estimator, where the Ridge estimator is obtained as a special case when an exchangeable prior around zero is assumed for the regression coefficients. Furthermore, it will be shown that the Efron-Morris (1972) estimators, at least for the known covariance cases, can be obtained as special cases of the estimator suggested here. Hence the latter can be regarded as an extension of Stein's method to the multivariate regression model.

¹ Brown and Zidek stack their matrices row by row. I find it more natural to stack them by columns and hence the few apparent differences.

The unknown covariance case, even after imposing initial restrictions on its structure, is shown to be nonidentifiable, and hence further (nonstochastic) restrictions are required for estimating the model parameters. Thus, prior information on at least some of the covariances is imposed in the form of a natural conjugate distribution.

2. THE LINEAR HIERARCHICAL MULTIVARIATE REGRESSION ESTIMATOR

The approach here is that of Lindley and Smith (1972). A three-stage model will be assumed, where in the first stage, given the design matrix X ,

$$y| \beta, \Gamma \sim N\left\{ \left(\begin{matrix} I \\ X \end{matrix} \right) \beta, \begin{pmatrix} \Gamma \\ \Gamma \otimes I_n \end{pmatrix} \right\}. \quad (5)$$

Equation (5) is Zellner's (1961) "seemingly unrelated regression equations" representation of the multivariate linear regression model in the special case of identical design matrices.

A second stage is added by supposing

$$\beta_j | \xi_j, \Omega \sim N\left(\begin{pmatrix} 1 \\ \xi_j \end{pmatrix}, \Omega \right), \quad \text{independent of } y, \quad (6)$$

that is, each $E\{\beta_j\} = \xi_j$ ($j=1, \dots, q$). A vague prior knowledge on ξ_j is assumed for the third stage.

The formal posterior distribution of $(\beta^T, \xi^T)^T$ given y and X is shown, for known variances, to be normal with mean vector [see Haitovsky (1985)] :

$$\begin{pmatrix} \beta^* \\ \xi^* \end{pmatrix} = \left(\begin{matrix} A^T C^{-1} A \\ A^T C^{-1} A \end{matrix} \right)^{-1} \begin{pmatrix} A^T C^{-1} y \\ A^T C^{-1} z \end{pmatrix} \quad (7)$$

and covariance matrix

$$\text{cov} \begin{pmatrix} \beta^* \\ \xi^* \end{pmatrix} = \left(\begin{matrix} A^T C^{-1} A \\ A^T C^{-1} A \end{matrix} \right)^{-1} \quad (8)$$

where
$$A_{\sim} = \begin{bmatrix} I_{\sim q} \otimes X_{\sim} & 0_{\sim} \\ -I_{\sim pq} & I_{\sim q} \otimes I_{\sim p} \end{bmatrix}, \quad C_{\sim} = \begin{pmatrix} \Gamma_{\sim} \otimes I_{\sim n} & 0_{\sim} \\ 0_{\sim} & \Omega_{\sim} \end{pmatrix},$$

and
$$z_{\sim}^T = (y_{\sim}^T, 0_{\sim}^T).$$

After some algebra this reduces to the means and variance, conditional on y_{\sim} and X_{\sim} :

$$\beta_{\sim}^* = \{ I_{\sim q} \otimes X_{\sim}^T X_{\sim} + (\Gamma_{\sim} \otimes I_{\sim p}) Q_{\sim} \}^{-1} (I_{\sim q} \otimes X_{\sim}^T) y_{\sim} \quad (9)$$

with

$$\text{cov}(\beta_{\sim}^*) = (\Gamma_{\sim}^{-1} \otimes X_{\sim}^T X_{\sim} + Q_{\sim})^{-1} \quad (10)$$

and

$$\xi_{\sim}^* = \left\{ (I_{\sim q} \otimes I_{\sim p})^T \Omega_{\sim}^{-1} (I_{\sim q} \otimes I_{\sim p}) \right\}^{-1} (I_{\sim q} \otimes I_{\sim p})^T \Omega_{\sim}^{-1} \beta_{\sim}^* \quad (11)$$

with

$$\text{cov}(\xi_{\sim}^*) = \left[(I_{\sim q}^T \otimes I_{\sim p}) \{ \Gamma_{\sim} \otimes (X_{\sim}^T X_{\sim})^{-1} + \Omega_{\sim}^{-1} \} (I_{\sim q} \otimes I_{\sim p}) \right]^{-1} \quad (12)$$

where

$$Q_{\sim} = [I_{\sim pq} - \Omega_{\sim}^{-1} (I_{\sim q} \otimes I_{\sim p}) \{ (I_{\sim q} \otimes I_{\sim p})^T \Omega_{\sim}^{-1} (I_{\sim q} \otimes I_{\sim p}) \}^{-1} (I_{\sim q} \otimes I_{\sim p})^T] \Omega_{\sim}^{-1}. \quad (13)$$

Equations (9) and (10) can be regarded as ridge estimators of the multivariate regression model. Furthermore, the MRRE reduces to the least squares estimator (3) when $\Omega_{\sim}^{-1} = 0$, as is the case with the univariate ridge.

It will be instructive at this stage to impose some structure on the general covariance matrix Ω_{\sim} . Let,

$$\Omega_{\sim} = V_{\sim} \otimes \Sigma_{\sim} \quad (14)$$

so that V_{\sim} ($q \times q$) > 0 is the covariance (except for an overall scale factor) between any two rows in B_{\sim} and Σ_{\sim} ($p \times p$) > 0 is the covariance (except for

an overall scale factor between any two columns in β . As a result Q reduces to $Q = V^{-1} W \otimes \Sigma^{-1}$, and hence eqns. (9) - (13) reduce respectively to

$$\beta^* = (I_{\nu q} \otimes X^T X + K_{\nu w} \otimes \Sigma^{-1})^{-1} (I_{\nu q} \otimes X^T) y \quad (15)$$

$$\text{cov}(\beta^*) = (I_{\nu q} \otimes X^T X + K_{\nu w} \otimes \Sigma^{-1})^{-1} (I_{\nu q} \otimes I_{\nu p}) \quad (16)$$

$$\xi^* = (1_{\nu q}^T V^{-1} 1_{\nu q})^{-1} (1_{\nu q}^T V^{-1} \otimes I_{\nu p}) \beta^* \quad \text{and} \quad (17)$$

$$\text{cov}(\xi^*) = [(1_{\nu q}^T \otimes I_{\nu p}) \{I_{\nu q} \otimes (X^T X)^{-1} + V \otimes \Sigma\}^{-1} (1_{\nu q} \otimes I_{\nu p})]^{-1}, \quad (18)$$

where $K_{\nu w} = \Gamma V^{-1} W$, $W = I_{\nu q} - (1_{\nu q}^T V^{-1} 1_{\nu q})^{-1} 1_{\nu q} 1_{\nu q}^T V^{-1}$ and $J_{\nu q} = 1_{\nu q} 1_{\nu q}^T$.

Brown & Zidek (1980) require that MRRE obeys two rules:

- (i) that it should be a Bayes rule for fixed known K_{ν} , which, being a posterior mean, our estimator (13) certainly is;
- (ii) for a suitably chosen estimator for K_{ν} , in the special case of equal information ($X^T X = I_{\nu}$), the MRRE should correspond to the Efron and Morris (1972) multivariate extension of the James and Stein (1961) estimator. It will be shown in the next section that this requirement is met here as well.

3. SPECIAL CASES

In this section we consider a few important subclasses of priors on the regression coefficients. We start with those affecting the ridge matrix $K_{\nu w}$.

3.1 The prior $\xi = 0$ is of special interest in applied work, since it results in the multivariate regression analog to the Hoerl and Kennard (1970) "ordinary" ridge rule, in Thisted's (1976) terminology. Thus, as Lindley and Smith (1972, p. 11) show, the univariate (multiple) regression estimator, when the elements of β_j are assumed exchangeable about an unknown value (i.e. a vague ξ), is $\beta_j^{LS} = (X_{\sim}^T X_{\sim} + k U_p)^{-1} X_{\sim}^T Y$, where $U_p = I_p - p^{-1} J_p$ and $k = \gamma/\sigma_{\beta}^2$, the ratio of the regression error variance to the prior variance of β . When exchangeability is assumed about 0 (i.e., $\xi = 0$, in our notation),

U_p reduces to I_p and the "ordinary" ridge is obtained. Similarly, in multivariate regression, $\xi = 0$ is implemented by deleting the last column of submatrices from A_{\sim} (eqn. 7), resulting in $Q = Q_{\sim}$ in eqn. (13). Or, when (14) is taken into account, in analogy to the univariate case, W_{\sim} becomes I_{\sim} and $K_{\sim w} = \Gamma V_{\sim}^{-1}$ (in contrast to $K_{\sim w} = \Gamma V_{\sim}^{-1} W_{\sim}$ when the prior on ξ is vague). $K_{\sim w}$ is an obvious multivariate generalization of k in the Lindley and Smith framework.

3.2 Multivariate Exchangeability Between Multiple Regression Equations is expressed by letting $V_{\sim} = I_{\sim}$, obtaining a ridge matrix $K_{\sim u} = \Gamma U_{\sim q}$.

If, furthermore, it is assumed that $\Gamma_{\sim} = \gamma I_{\sim q}$ (15) reduces to

$$\beta_{\sim i}^* = H_{\sim i} \hat{\beta}_{\sim i} + (I_{\sim p} - H_{\sim}) q^{-1} \sum_{j=1}^q \hat{\beta}_{\sim j}$$

with

$$\text{cov}(\beta_{\sim i}^*) = \{H_{\sim} + q^{-1}(I_{\sim p} - H_{\sim})\} \gamma (X_{\sim}^T X_{\sim})^{-1}$$

where $H_{\sim} = (\gamma^{-1} X_{\sim}^T X_{\sim} + \Sigma_{\sim}^{-1})^{-1} \gamma^{-1} X_{\sim}^T X_{\sim}$ and $\hat{\beta}_{\sim i}$ as before are the least squares estimators. This is the Bayesian approach to growth curves with identical X_i matrices as developed by Fearn [1975, equations (10) and (11)].

It is worth mentioning here that if Zellner's seemingly unrelated regression

equations model is assumed instead of (5), i.e., the design matrix Z_{\sim} , say, in (5) takes the form of a diagonal block matrix with different $X_{\sim i}$ ($i=1, \dots, q$) submatrices along the main diagonal and exchangeability between the regression equations is assumed, then

$$\beta_{\sim}^* = \{Z_{\sim}^T (\Gamma_{\sim}^{-1} \otimes I_{\sim n}) Z_{\sim} + U_{\sim q} \otimes \Sigma_{\sim}^{-1}\}^{-1} Z_{\sim}^T (\Gamma_{\sim}^{-1} \otimes I_{\sim n}) y_{\sim}.$$

See Lindley and Smith (1972, eq. 23) and Haitovsky (1979, §VI-1).

Here again if Γ_{\sim} is a diagonal matrix, Fearn's results for non-identical design matrices are easily obtained. See also Lindley and Smith (1972, p.10).

3.3 Multivariate Exchangeability Within Multiple Regression Equations will require, in addition to letting $\Sigma_{\sim} = I_{\sim p}$, a fourth stage in the Linear Hierarchical Model [cf. Haitovsky (1979)]. If only the former is assumed then,

$$\beta_{\sim}^* = (I_{\sim q} \otimes X_{\sim}^T X_{\sim} + K W \otimes I_{\sim p})^{-1} (I_{\sim q} \otimes X_{\sim}^T) y_{\sim} \quad (19)$$

The last two cases can be combined with the $\xi_{\sim} = 0$ prior assumption, resulting in $W_{\sim} = U_{\sim q} = I_{\sim q}$ and hence, for instance for the present case one obtains:

$$\beta_{\sim}^* = (I_{\sim q} \otimes X_{\sim}^T X_{\sim} + K \otimes I_{\sim p})^{-1} (I_{\sim q} \otimes X_{\sim}^T) y_{\sim} = \hat{\beta}_{\sim}^*(K). \quad (20)$$

This is the Brown-Zidek MRRE.

3.4 The Efron-Morris Extension of Stein's Method can be reduced from the hierarchical representation of the multivariate regression model (5) and (6) by letting $n=p$ and $X_{\sim} = I_{\sim p}$, $\Gamma_{\sim} = I_{\sim q}$, $\xi_{\sim} = 0$ and $\Omega_{\sim} = I_{\sim q} \otimes \Sigma_{\sim}$, thus obtaining

$$\beta_{\sim}^* = (I_{\sim pq} + I_{\sim q} \otimes \Sigma_{\sim}^{-1})^{-1} y_{\sim} = \{I_{\sim pq} - (I_{\sim pq} + I_{\sim q} \otimes \Sigma_{\sim}^{-1})^{-1} (I_{\sim q} \otimes \Sigma_{\sim}^{-1})\} y_{\sim}$$

or

$$\beta_{\sim i}^* = \{I_{\sim p} - (\Sigma_{\sim} + I_{\sim p})^{-1}\} y_{\sim i} \quad (i=1, \dots, q).$$

Cf., Efron and Morris (1972, eqn. 2.4).

If we relax now the *a priori* assumption that $\xi = 0$ and let it be unknown [Efron and Morris (1972, §7)], ξ is estimated according to eqn. (17)

by $\bar{y} = \frac{1}{p} \sum_{i=1}^p y_i$ and hence

$$\beta_{\lambda}^* = (I_{\lambda pq} + I_{\lambda q} \otimes \Sigma_{\lambda}^{-1})^{-1} (y + \frac{1}{\lambda q} \otimes \Sigma_{\lambda}^{-1} \bar{y})$$

or

$$\begin{aligned} \beta_{\lambda i}^* &= (I_{\lambda p} + \Sigma_{\lambda}^{-1})^{-1} (y_{\lambda i} + \Sigma_{\lambda}^{-1} \bar{y}) = (I_{\lambda p} + \Sigma_{\lambda}^{-1})^{-1} \{y_{\lambda i} - \bar{y} + (I_{\lambda p} + \Sigma_{\lambda}^{-1}) \bar{y}\} \\ &= \{I_{\lambda p} - (I_{\lambda p} + \Sigma_{\lambda}^{-1})^{-1}\} (y_{\lambda i} - \bar{y}) + \bar{y} \quad (i=1, \dots, q) \end{aligned}$$

Cf. Efron and Morris (1972, eqn. 7.2) and recall that $(k-p-2)^{-1} \sum_{\lambda}^{\infty}$ there estimates $I_{\lambda p} + \Sigma_{\lambda}$.

Efron and Morris's more general formulation (1972, §4) requires a modification in our model, i.e., replacing the covariance matrix in (5) by $I_{\lambda q} \otimes \Gamma_{\lambda}$ and setting $X = I_{\lambda p}$, $V = I_{\lambda q}$, one obtains

$$\beta_{\lambda}^* = \{I_{\lambda q} \otimes (\Gamma_{\lambda}^{-1} + \Sigma_{\lambda}^{-1})^{-1} \Gamma_{\lambda}^{-1}\} y_{\lambda}$$

or

$$\beta_{\lambda i}^* = \{I_{\lambda p} - \Gamma_{\lambda} (\Gamma_{\lambda} + \Sigma_{\lambda})^{-1}\} y_{\lambda i} \quad (i=1, \dots, q)$$

for the $\xi = 0$ case [cf. Efron and Morris (1972, eqn. 4.6)], and

$$\beta_{\lambda i}^* = \bar{y} + \{I_{\lambda p} - \Gamma_{\lambda} (\Gamma_{\lambda} + \Sigma_{\lambda})^{-1}\} (y_{\lambda i} - \bar{y}) \quad (i=1, \dots, q)$$

for the unknown ξ case.

4. RESTRICTED LEAST SQUARES INTERPRETATION

Ridge regression estimators were originally obtained by either minimizing the residual error sum of squares subject to a positive (scalar) constraint on the length of the vector of regression parameters, or by minimizing the length of the parameter vector subject to a positive constraint on the residual error sum of squares. A generalization is possible if the minimization of the residual error sum of squares is performed subject to a constraint on the length, in a prespecified positive metric, of the *differences* between the vector of parameters and an *a priori* given vector of constants. See Haitovsky and Wax (1980). Similarly, multivariate ridge-regression-like estimates can be obtained as a solution to the problem:

$$\min_{\beta} \{ y - (I_q \otimes X) \beta \}^T (\Gamma \otimes I_n)^{-1} \{ y - (I_q \otimes X) \beta \}$$

subject to the constraint

$$(\beta - 1_q \otimes \bar{\beta})^T \Omega^{-1} (\beta - 1_q \otimes \bar{\beta}) = h$$

yielding

$$\hat{\beta}^*(k, \Omega^{-1}, \bar{\beta}) = (\Gamma^{-1} \otimes X^T X + k \Omega^{-1})^{-1} \{ (\Gamma^{-1} \otimes X^T) y + k \Omega^{-1} (1_q \otimes \bar{\beta}) \}$$

where $\bar{\beta}$ is a $p \times 1$ vector of constants, h is a given positive constant, and $k > 0$ is the Lagrange multiplier. If now instead of fixing $\bar{\beta}$ one takes it to be (in some sense) a mean vector of the estimated β_j 's, one obtains (9) with Q replaced by kQ , where Q depends on how exactly $\bar{\beta}$ was derived. If, on the other hand, $\bar{\beta}$ is set to 0

$$\hat{\beta}^*(k, \Omega^{-1}, 0) = (\Gamma^{-1} \otimes X^T X + k \Omega^{-1})^{-1} (\Gamma^{-1} \otimes X^T) y,$$

i.e., (9) and (10) are obtained with Q replaced by $k \Omega^{-1}$.

In a similar manner the ridge regression estimators are obtained from the generalized ridge estimators by setting the *a priori* given vector of constants equal to zero.

The properties of the generalized ridge regression estimators are explored by Haitovsky and Wax (1980).

5. COVARIANCES UNKNOWN

Turning now to the more realistic situation where the covariances are unknown, it is natural to choose maximum likelihood as the estimation principle. One may regard our three-stage hierarchical model as a mixed effect model, where the first stage parameters are assumed random and the second stage parameters are assumed fixed, in view of the widely known fact that classical (BLUE) estimators of fixed effects are the limiting posterior means corresponding to random effects with zero precision matrix (i.e. flat prior). Maximum likelihood estimators for the mixed effect models were reviewed by Harville (1977) for variance component models and were extended by Dempster, Rubin and Tsutakawa (1979) for covariance component models. We follow Haitovsky (1985) who writes the joint probability function, when assuming (14), as proportional to

$$|\Gamma_{\sim} \otimes I_{\sim}|^{-1/2} |V_{\sim} \otimes \Sigma_{\sim}|^{-1/2} \text{etr} \left\{ -\frac{1}{2} (Y_{\sim} - XB_{\sim}) \Gamma_{\sim}^{-1} (Y_{\sim} - XB_{\sim})^T - \frac{1}{2} \Sigma_{\sim}^{-1} \Xi_{\sim} V_{\sim}^{-1} \Xi_{\sim}^T \right\} \quad (21)$$

where etr denotes the exponent of the matrix trace, and $\Xi_{\sim} = B_{\sim} - \frac{1}{\lambda q} \otimes \xi_{\sim}$ is obtained by "rematricizing" the $\beta_{\sim} - \frac{1}{\lambda q} \otimes \xi_{\sim}$ vector. Maximizing (21) with respect to Γ_{\sim}^{-1} , V_{\sim}^{-1} , Σ_{\sim}^{-1} , B_{\sim} and ξ_{\sim} , the following estimators are obtained:

$$\left. \begin{aligned} \hat{\Gamma}_{\sim} &= \frac{1}{n} (Y_{\sim} - XB_{\sim}^*)^T (Y_{\sim} - XB_{\sim}^*) \\ \hat{V}_{\sim} &= \frac{1}{p} \hat{\Xi}_{\sim}^T \hat{\Sigma}_{\sim}^{-1} \hat{\Xi}_{\sim} \\ \hat{\Sigma}_{\sim} &= \frac{1}{q} \hat{\Xi}_{\sim} \hat{V}_{\sim}^{-1} \hat{\Xi}_{\sim}^T \end{aligned} \right\} \quad (22)$$

and $\hat{\beta}_{\sim}^*$ and $\hat{\xi}_{\sim}^*$ are as in (15) and (17), but with the covariance matrices

replaced by their estimators.

There are two things to note here:

- a) in the traditional treatment of the multivariate regression, the matrix of coefficients B_{\sim} is considered fixed, hence V_{\sim} and Σ_{\sim} are equal to zero, and Γ_{\sim} is indeed estimated by $\hat{\Gamma}_{\sim}$ in (22). Likewise, if it is assumed that $\Sigma_{\sim} = I_{\sim p}$ and $\xi_{\sim} = 0$, i.e., the Brown-Zidek MRRE (see Section 3.3 above), there is no need to estimate Σ_{\sim} , while V_{\sim} is estimated by

$$\hat{V}_{\sim} = \frac{1}{p} \hat{B}_{\sim}^* T \hat{B}_{\sim}^* \quad (23)$$

where \hat{B}_{\sim}^* is a matrix of estimators obtained by "rematricizing" $\hat{\beta}_{\sim}^*$ with $\xi_{\sim} = 0$ and $\Sigma_{\sim} = I_{\sim p}$. Estimator (23) is the "natural estimator", suggested by Brown and Zidek (1980, eqn. 2.14), apart from some differences in the number of degrees of freedom, which the ML method so notoriously fails to account for. A correction for the number of degrees of freedom is suggested below which will make (23) identical to the Brown-Zidek estimator.

- b) the matrix of order $p \times q$

$$\hat{\Xi}_{\sim} = \hat{B}_{\sim}^* - \frac{1}{q} \hat{B}_{\sim}^* \otimes \xi_{\sim} = \hat{B}_{\sim}^* \{ I_{\sim q} - (\frac{1}{q} \hat{V}_{\sim}^{-1} \frac{1}{q})^{-1} \hat{V}_{\sim}^{-1} J_{\sim} \}$$

is of rank $\leq \min \{p, q-1\}$, since its rows lie in the $q-1$ dimensional space. (Note that $W_{\sim}^T = I_{\sim q} - (\frac{1}{q} \hat{V}_{\sim}^{-1} \frac{1}{q})^{-1} \hat{V}_{\sim}^{-1} J_{\sim}$ is an idempotent (nonsymmetric) matrix with rank and trace equal to $q-1$). Consequently, either $\hat{\Sigma}_{\sim}$ or \hat{V}_{\sim} is a singular matrix. Alternatively, there are $p(p+1)/2 + q(q+1)/2 - 1$ independent parameters to be estimated in $V_{\sim} \otimes \Sigma_{\sim}$, but only pq ($< p(p+1)/2 + q(q+1)/2 - 1$) entries in \hat{B}_{\sim}^* . That is, the structure imposed on Ω_{\sim} does not suffice to identify the second stage covariance matrix.

In order to overcome the singularity, a Bayesian solution is suggested.

To this end we specify prior distributions for the covariance matrices. It will be done through the conjugate family, which is the Wishart distribution. Assuming Γ_{\sim}^{-1} , V_{\sim}^{-1} and Σ_{\sim}^{-1} are independent, Wishart distributed with γ , ν and σ degrees of freedom and G_{\sim} , R_{\sim} and S_{\sim} matrices, respectively, the following densities (up to a proportionality factor) are obtained:

$$|\Gamma_{\sim}|^{-\frac{1}{2}(\gamma-q-1)} \text{etr}\{-\frac{1}{2}\Gamma_{\sim}^{-1}G_{\sim}\}$$

$$|V_{\sim}|^{-\frac{1}{2}(\nu-q-1)} \text{etr}\{-\frac{1}{2}V_{\sim}^{-1}R_{\sim}\}$$

and

$$|\Sigma_{\sim}|^{-\frac{1}{2}(\sigma-p-1)} \text{etr}\{-\frac{1}{2}\Sigma_{\sim}^{-1}S_{\sim}\}.$$

The joint distribution of all the quantities is obtained by multiplying the joint probability function (21) by the three prior densities above. The integration with respect to ξ_{\sim} is performed by rewriting the second term in the exponent in (21) as

$$\begin{aligned} \text{tr}\{\Sigma_{\sim}^{-1}EV_{\sim}^{-1}\bar{\xi}_{\sim}^T\} &= (\beta_{\sim} - \frac{1}{\sim q} \otimes \xi_{\sim})^T (V_{\sim} \otimes \Sigma_{\sim})^{-1} (\beta_{\sim} - \frac{1}{\sim q} \otimes \xi_{\sim}) \\ &= (\beta_{\sim} - \frac{1}{\sim q} \otimes \bar{\xi}_{\sim})^T (V_{\sim} \otimes \Sigma_{\sim})^{-1} (\beta_{\sim} - \frac{1}{\sim q} \otimes \bar{\xi}_{\sim}) + (\frac{1}{\sim q} V_{\sim}^{-1} \frac{1}{\sim q}) (\xi_{\sim} - \bar{\xi}_{\sim})^T \Sigma_{\sim}^{-1} (\xi_{\sim} - \bar{\xi}_{\sim}) \end{aligned}$$

where $\bar{\xi}_{\sim} = (\frac{1}{\sim q} V_{\sim}^{-1} \frac{1}{\sim q})^{-1} (\frac{1}{\sim q} V_{\sim}^{-1} \otimes I_{\sim p}) \beta_{\sim}$ [the counterpart of (17)], since the cross product term vanishes.

Now the integration with respect to ξ_{\sim} is straightforward, resulting in the joint posterior density for β_{\sim} , Γ_{\sim} , V_{\sim} and Σ_{\sim} which is proportional to:

$$\begin{aligned} &|\Gamma_{\sim}|^{-\frac{1}{2}(n+\gamma-q-1)} |V_{\sim}|^{-\frac{1}{2}(p+\nu-q-1)} |\Sigma_{\sim}|^{-\frac{1}{2}(q+\sigma-p-2)} (\frac{1}{\sim q} V_{\sim}^{-1} \frac{1}{\sim q})^{-\frac{1}{2}} \\ &\times \text{etr}\{-\frac{1}{2}\{\Gamma_{\sim}^{-1}G_{\sim} + V_{\sim}^{-1}R_{\sim} + \Sigma_{\sim}^{-1}S_{\sim} + (Y_{\sim} - XB)_{\sim} \Gamma_{\sim}^{-1} (Y_{\sim} - XB)_{\sim}^T + \Sigma_{\sim}^{-1}EV_{\sim}^{-1}\bar{\xi}_{\sim}^T\}\} \end{aligned} \quad (24)$$

where here $\bar{\xi} = B - \frac{1}{n} \otimes \bar{\xi}$. The integration with respect to ξ resulted in the loss of one degree of freedom in the estimation of ξ .

Following Lindley and Smith (1972) we derive the modal Bayesian estimators by differentiating the log joint posterior density partially with respect to B , Γ^{-1} , V^{-1} and Σ^{-1} obtaining \hat{B}^* as before (eqn. 15) with Γ , Σ and V replaced by their estimators:

$$\hat{\Gamma} = (n + \gamma - q - 1)^{-1} \{ G + (Y - XB^*)^T (Y - XB^*) \}$$

$$\hat{\Sigma} = (q + \sigma - p - 2)^{-1} \{ S + \hat{\Xi} \hat{V}^{-1} \hat{\Xi}^T \}$$

$$\hat{V} = (p + \nu - q - 1)^{-1} \{ (\frac{1}{n} \otimes \hat{V}^{-1} \frac{1}{n})^{-1} J_{nq} + R + \hat{\Xi} \hat{\Sigma}^{-1} \hat{\Xi} \}$$

$$= (p + \nu - q - 1)^{-1} [\{ (p + \nu - q - 2) \frac{1}{n} \otimes (R + \hat{\Xi} \hat{\Sigma}^{-1} \hat{\Xi})^{-1} \frac{1}{n} \}^{-1} J_{nq} + R + \hat{\Xi} \hat{\Sigma}^{-1} \hat{\Xi}]$$

The estimators for ξ and V are natural in view of the following theorem:

Theorem: Let the $p \times q$ random matrix B be multivariate normal with a common mean vector ξ and variance $V \otimes \Sigma$, $V, \Sigma > 0$, i.e., $\beta = \text{vec}\{B\} \sim N(\frac{1}{n} \otimes \xi, V \otimes \Sigma)$. Then, if $q \geq p + 1$

$$\hat{\Xi} \hat{V}^{-1} \hat{\Xi}^T \sim W(\Sigma, p, q);$$

and if $p \geq q + 1$, then

$$\hat{\Xi}^T \hat{\Sigma}^{-1} \hat{\Xi} \sim W(V, q, p)$$

where here $\Xi = B - \frac{1}{n} \otimes \xi$ is the matrix of columnwise deviations from the common mean factor.

Proof: $\text{vec } \{\Xi\} \sim N(0, V \otimes \Sigma)$. Diagonalize V by the orthogonal matrix P , i.e., $P^T V P = \Lambda$, a diagonal matrix, and define $z = (\Lambda^{-1/2} \otimes I_p)(P^T \otimes I_p) \text{vec } \{\Xi\} = (\Lambda^{-1/2} P^T \otimes I_p) \text{vec } \{\Xi\}$. Clearly, $z \sim N(0, I_q \otimes \Sigma)$, hence $ZZ^T \sim W(\Sigma, p, q)$, where $z = \text{vec } \{Z\}$, but $ZZ^T = EP^T \Lambda^{-1} P \Xi^T = EV^{-1} \Xi^T$. Similarly for $p \geq q+1$.

Q.E.D.

Corollary: For $q \geq p+1$, $(\Xi V^{-1} \Xi^T)^{-1} \sim W^{-1}(\Sigma^{-1}, p, p+q+1)$ and hence, $(q-p-1) E\{(\Xi V^{-1} \Xi^T)^{-1}\} = \Sigma^{-1}$. If $p \geq q+1$, then $(\Xi \Sigma^{-1} \Xi^T)^{-1} \sim W^{-1}(V^{-1}, q, p+q+1)$ and hence $(p-q-1) E\{(\Xi \Sigma^{-1} \Xi^T)^{-1}\} = V^{-1}$.

The correction of the number of degrees of freedom in the estimation of Γ differs from that in least squares theory; in the latter the divisor is $(n-p-1)$.

This brings up a related problem. O'Hagan (1976) proves that in joint posterior densities as in our (24), modal estimators obtained from the joint density will differ, often substantially, from that of the marginal densities; the main difference being the number of degrees of freedom, i.e., the divisors of the estimated variances and covariances. Thus, the problem of which mode to choose can be problematic, especially when the number of regression coefficients to be estimated is large relative to the sample size. However, note that this issue affects the estimation of the vector of regression coefficients only indirectly, via the estimation of the variances and covariances. It is our opinion that the question of "which mode" cannot be ignored, nor overstated in situations like ours, where the main objective is the estimation of the regression coefficients rather than the covariances. A few trials made with our model and the example in O'Hagan (1976) indeed support this contention.

The Bayesian modal estimates can be calculated using an iterative procedure: initial values for \hat{B}_0^* , \hat{B}_1^* , are computed using least squares multivariate regression estimates (eqn. 3), from which the initial values for $\hat{\Gamma}$ and for $\hat{\xi}^*$ (eqn. 17) are computed, using for the latter a simple average of the q columns of \hat{B}_0^* (i.e., letting $\hat{V} = I_q$ in eqn. 17). Next, $\hat{\Xi}_0$ is computed, from which

$\hat{\Sigma}_{\lambda_0}$ is obtained after suppressing \hat{Y}_{λ} to I_{λ_q} . Having computed $\hat{\Sigma}_{\lambda_0}$ and $\hat{\lambda}_0$, \hat{Y}_{λ_0} can be calculated, completing the first round, from which iterations can start.

Another possible iterative solution is the E-M algorithm. See Dempster et al. (1979)

6. APPLICATION TO SCOTTISH ELECTION.

Merely as an illustration, our method is applied to the Scottish election of October 1974 in forecasting the final results on the basis of partial returns. The forecasting model is that of Brown and Zidek (1979), which was used as purely illustrative as well. It is a modified version of a model actually used by Brown and Payne (1975) for predicting the February 1974 British General election.

6.1 Description of the Data: The data are taken from Brown and Zidek (1979) and are described there in detail. The forecasting problem is to forecast the winner of the undeclared constituencies on the basis of the declared constituencies.

The forecasting model is a multivariate regression with the dependent variables being the change in the party's share in the number of electorates in each constituency, relative to previous election. The independent variables are the party's shares of votes in the previous election plus three dummy variables. The four main parties who participated in the February and October 1974 elections were the Conservative, Labour, Liberal, and Nationalist parties. They will be denoted by $i = 1, 2, 3, 4$ in that order. Thus, let W_i and Z_i denote the October 1974 and February 1974 votes for the i^{th} party, respectively, and E the electorate figure, which stayed practically the same in the two elections. Then

$$X_i = Z_i/E \quad \text{and} \quad Y_i = W_i/E - X_i \quad i=1, \dots, 4$$

denote the (non dummy) explanatory variables and the dependent variables, respectively.

The 3 dummy variables are

$$X_5 = \begin{cases} 0.5 & \text{Liberal intervenes, i.e. } w_3 > 0 \text{ and } X_3 = 0; \\ 0 & \text{otherwise;} \end{cases}$$

$$X_6 = \begin{cases} 0.5 & R = 5, 6; \\ 0 & \text{otherwise;} \end{cases}$$

$$X_7 = \begin{cases} 0.5 & \text{Labour or Nationalist top party in February 1974} \\ & \text{and } |X_2 - X_4| \leq 0.2, \\ 0 & \text{otherwise.} \end{cases}$$

where R is a categorical variable defining region where

- | | |
|----------------|------------------------------------|
| 1 = GLASGOW; | 2 = Rest of Clydeside conurbation; |
| 3 = EDINBURGH; | 4 = Rest of industrial centres; |
| 5 = HIGHLANDS; | 6 = Rest of Scotland. |

The value of 0.5 employed in these three dummy variables X_5, X_6, X_7 is somewhat arbitrary but was chosen so that a priori coefficients for all seven variables would be of a similar magnitude.

Thus, $q = 4$ and $p = 8$, including the unit vector for the intercepts.

6.2 The Choice of Priors in the multivariate regression with a covariance structure (14) is being studied now and will be reported in a subsequent paper. The choice of the priors in the present paper motivated the research and, hence is partially based on it.

There are three covariance matrices G_{\sim} , R_{\sim} , and S_{\sim} and three scalars γ , ν , and σ , which require prior specification. Vague priors are assumed for the regression error covariance, i.e., G_{\sim} and γ are set to zero. The remaining prior values are related to the regression coefficients and their assessment should be based on their careful interpretation:

The forecasting equations discussed in the last section can be viewed as including an implied explanatory variable for the eligible voters who did not participate in the February election. Thus, rewriting slightly the first equation as:

$$\begin{aligned} \frac{W_{1j} - Z_{1j}}{E} = & \beta_{01} + \beta_{11} \frac{Z_{1j}}{E} + \beta_{21} \frac{Z_{2j}}{E} + \beta_{31} \frac{Z_{3j}}{E} + \beta_{41} \frac{Z_{4j}}{E} \\ & + \gamma_1 \frac{E - Z_{1j} - Z_{2j} - Z_{3j} - Z_{4j}}{E} + \text{dummies} + \text{error} \quad (j=1, \dots, 71) \end{aligned}$$

or, collecting terms:

$$\begin{aligned} \frac{W_{1j}}{E} = & (\beta_{01} + \gamma_1) + (1 + \beta_{11} - \gamma_1) \frac{Z_{1j}}{E} + (\beta_{21} - \gamma_1) \frac{Z_{2j}}{E} + (\beta_{31} - \gamma_1) \frac{Z_{3j}}{E} \\ & + (\beta_{41} - \gamma_1) \frac{Z_{4j}}{E} + \text{dummies} + \text{error} \end{aligned}$$

and similarly for the three other parties, it can be seen that the coefficients measure the shifts of votes from last election to the present, adjusted for the fraction of abstainer, attributable to the various party supporters in the last election. Thus, for instance, the negative β_4 in the "Conservative equation" signifies a shift of those who voted Nationalist in the last election away from the Conservative party in the present election, over and above the shift away from the conservative party of nonvoters in the last election. If we disregard now the problem of abstainers, obviously each of the regression coefficients across equations must sum up to zero: in the absence of nonvoters the reallocation of last election voters for each party will result in compensating shifts in the present election. In the absence of more specific information the componentwise deviation of the β -vectors around ξ is justified in the present application. Moreover, $\xi = 0$ should be imposed if all eligible voters vote, or if it is believed that there will be no (non-proportional) shifts among the last election abstainers, or nearly so. Otherwise, let ξ be freely estimated.

Now, as the case happened to be, competition was particularly fierce between the Conservative and Nationalist parties, but less so between the Labour

and Nationalist. In the final result the Conservatives lost in the October election 28.35% of their voters in February, the Labour lost 5.39%, the Nationalist gained 32.61%, while the Liberals remained the same. The participation rate dropped from 77.87% in February to 74.60% in October.

Assuming that the nature of competition between the Conservatives and the Nationalists on the one hand, and the Labour and the Nationalists on the other, was known before the election, we have chosen a priori values for R_{\sim} to reflect a prior belief that $\text{corr}(\beta_{i1}, \beta_{i4}) = -0.7$ and $\text{corr}(\beta_{i2}, \beta_{i4}) = -0.5$; all other correlations between the β 's across equations were set to zero. Furthermore, in order to be in line with $\Xi_{\sim}^T \Sigma_{\sim}^{-1} \Xi_{\sim}$, to which R_{\sim} is added, the diagonal elements of R_{\sim} were set to 0.1. Thus, $R_{\sim} = 0.1 I_{\sim 4} + R^*$, where R^* is a 4×4 matrix of zeros except that -0.07 appears in the (1,4) and (4,1) positions and -0.05 appears in the (2,4) and (4,2) positions.

As for the choice of S_{\sim} , it was felt that the Conservatives benefited from the Liberal intervention, but were adversely affected by the last two dummy variables, and hence we set $S_{\sim} = I_{\sim 8} + S^*$ where S^* is an 8×8 matrix of zeros except for $\frac{1}{2}$ in the (2,6) and (6,2) positions, $-\frac{1}{3}$ in the (2,7) and (7,2) positions, and $-\frac{1}{2}$ in the (2,8) and (8,2) positions. In the six explanatory variable case the non-zero off-diagonal elements were slightly increased to 0.7 and -0.5 respectively, while the last row and column were deleted. Finally, $v = 10$ and $\sigma = 20$ produced divisors for the estimates of V_{\sim} and Σ_{\sim} of 13 and 14 respectively, and judged to reflect our confidence in our choice of R_{\sim} and S_{\sim} .

6.3 The Forecasting Performance Criteria: It is not obvious that forecasting the winner is best achieved by first forecasting the vote shares from which the winner can be picked out. However, there are countries with proportional representation where the *number* of votes matters, and, hence the multivariate

regression might be better suited there. We nevertheless follow here, merely as an illustration, Brown and Zidek (1979) and use the first 25 declared constituencies to forecast the results of the remaining $71 - 25 = 46$ constituencies, using the multivariate regression set-up with n , number of observations, equal to 25. Using these 25 observations we first estimate $\hat{\beta}_i^*$ from which we can predict the Y_i 's and eventually the number of votes to be cast in each constituencies for the four parties. Hence,

$$\hat{W}_{ij} = (\hat{Y}_{ij} + X_{ij})E_i \quad i = 26, \dots, 71; \quad j = 1, \dots, 4.$$

Three criteria for goodness of prediction were chosen.

$$SD = \left\{ \text{tr}(\hat{W} - \hat{\hat{W}}) (\hat{W} - \hat{\hat{W}})^T / (4 \times 46) \right\}^{1/2}$$

PRED = number of incorrect predictions of winning party

$$RSD = \left\{ \sum_i \sum_j (W_{ij} - \hat{W}_{ij})^2 / (W_{ij}^2 \times 4 \times 46) \right\}^{1/2}$$

where \hat{W} is the 46×4 matrix with typical element \hat{W}_{ij} .

The first two measures were used by Brown and Zidek (1979). The first is just the square root of the mean squared prediction error and it indeed estimates the standard deviation of prediction.

The second records the number of times the prediction of winners proved to be wrong. It is a very crude measure, but has the appeal of simplicity. It does not however take account of the closeness of a particular contest.

The last is the square root of the *relative* mean squared prediction error. It standardizes the prediction errors, in some sense, for ease of prediction. Brown and Zidek also used a third measure which we feel will a priori bias the results in favour of their method, since it was constructed to be in line with the loss function they used in developing their estimates, and thus is rejected here in favour of more objective measures.

6.4 Results: Two variants of the linear hierarchical model were tried and compared to the multivariate regression least-squares (maximum likelihood) estimates (eqn. 3) and to the Brown-Zidek ridge estimates (eqn. 4 with κ estimated by $\hat{\Gamma}^{-1}\hat{Y}$, where $\hat{\Gamma}$ and \hat{Y} are given on p. 13 below except that in the latter the divisor is $p-q-1$). The two variants are with ξ vague and with ξ a prior set to zero. The three estimates are reported in Tables I - III. The predictive record of all those three estimates and that of Brown-Zidek are summarized in Table IV.

TABLE IV
PREDICTION RECORDS OF THE DIFFERENT ESTIMATES

	SD	RSD	PRED*
Least Squares (M.L.)	1551.03	0.2972	7(0)
Ridge (Brown-Zidek)	1315.78	0.2189	6(1)
Linear Hierarchical Model			
$\xi = 0$	1315.62	0.2271	6(1)
ξ vague	1423.97	0.2651	5(1)

* The numbers in parentheses denote the number of wrong "predictions" of the winner party in the 25 declared constituencies.

Interestingly, the best prediction records is for the Ridge and the linear hierarchical model with $\xi = 0$, which is closely related to Brown-Zidek's Ridge, since the latter essentially presupposes exchangeability about zero. See section 3. The reason for the more restrictive cases, namely $\xi = 0$, to perform better might be due to the fact that the "sample" of the 25 declared

TABLE I

LEAST SQUARES (ML) ESTIMATES OF THE MATRIX OF COEFFICIENTS OF THE MULTIVARIATE
REGRESSION MODEL AND THEIR RATIO TO THE CORRESPONDING STANDARD DEVIATIONS

	Equation 1		Equation 2		Equation 3		Equation 4	
	$\hat{\beta}$	$\hat{\beta}/\text{Std. Dev.}$	$\hat{\beta}$	$\hat{\beta}/\text{Std. Dev.}$	$\hat{\beta}$	$\hat{\beta}/\text{Std. Dev.}$	$\hat{\beta}$	$\hat{\beta}/\text{Std. Dev.}$
β_0	-0.0119	0.22	-0.0007	0.01	0.0367	0.69	-0.1501	2.06
β_1	-0.0627	1.04	-0.0911	2.32	-0.0303	0.51	0.1793	2.21
β_2	-0.0297	0.32	-0.0212	0.36	-0.0712	0.78	0.3193	2.54
β_3	-0.1070	0.79	0.2768	3.13	-0.2864	2.15	0.4736	2.59
β_4	-0.2202	1.88	-0.0768	1.01	-0.0854	0.74	0.3321	2.10
β_5	-0.0160	0.63	0.0276	1.67	0.0763	3.06	-0.0250	0.73
β_6	0.0315	0.92	0.0307	1.37	0.0028	0.08	-0.0146	0.31
β_7	-0.0061	0.36	-0.0025	0.22	0.0042	0.25	-0.0278	1.20

TABLE II

LINEAR HIERARCHICAL MODEL ESTIMATES OF THE MATRIX OF COEFFICIENTS OF THE
MULTIVARIATE REGRESSION MODEL AND THEIR RATIO TO THE CORRESPONDING STANDARD DEVIATIONS

	Equation 1		Equation 2		Equation 3		Equation 4		ξ	
	β^*	$\frac{\beta^*}{\text{Std. Dev.}}$	β^*	$\frac{\beta^*}{\text{Std. Dev.}}$	β^*	$\frac{\beta^*}{\text{Std. Dev.}}$	β^*	$\frac{\beta^*}{\text{Std. Dev.}}$	ξ^*	$\frac{\xi^*}{\text{Std. Dev.}}$
β_0	-0.0768	3.75	-0.0234	1.22	-0.0486	2.33	0.0364	1.56	-0.0183	0.47
β_1	-0.0254	0.97	-0.0338	1.45	0.0301	1.13	-0.0124	0.38	-0.0198	0.28
β_2	0.0452	1.18	0.0274	0.79	0.0619	1.59	0.0606	1.33	0.0477	0.94
β_3	0.0861	1.45	0.1858	3.48	-0.0450	0.74	0.1966	2.73	0.1474	2.87
β_4	-0.0721	1.52	-0.0533	1.22	-0.0157	0.32	-0.0300	0.54	-0.0485	0.84
β_5	-0.0022	0.16	0.0025	0.21	0.1024	7.00	-0.0184	0.95	-0.0026	0.05
β_6	0.0135	0.69	0.0330	2.09	0.0162	0.81	0.0136	0.51	0.0184	0.34

TABLE III

LINEAR HIERARCHICAL ESTIMATES OF THE MATRIX OF COEFFICIENTS OF THE MULTIVARIATE
REGRESSION AND THEIR RATIOS TO THEIR CORRESPONDING STANDARD DEVIATIONS

$$\xi = 0$$

	Equation 1		Equation 2		Equation 3		Equation 4	
	β^*	$\beta^*/\text{Std. Dev.}$	β^*	$\beta^*/\text{Std. Dev.}$	β^*	$\beta^*/\text{Std. Dev.}$	β^*	$\beta^*/\text{Std. Dev.}$
β_0	-0.0595	5.73	-0.0036	0.41	-0.0346	3.06	0.0574	4.99
β_1	-0.0166	0.92	-0.0208	1.35	0.0254	1.30	0.0083	0.41
β_2	0.0022	0.09	-0.0140	0.71	0.0218	0.87	0.0034	0.14
β_3	-0.0038	0.14	0.0236	1.01	-0.0538	1.79	0.0203	0.73
β_4	-0.0264	0.99	-0.0096	0.42	0.0083	0.28	0.0135	0.50
β_5	-0.0232	1.96	-0.0247	2.49	0.0935	7.60	-0.0345	2.17
β_6	-0.0075	0.47	0.0103	0.78	-0.0011	0.07	-0.0047	0.25
β_7	-0.0112	0.82	-0.0000	0.00	0.0084	0.59	-0.0135	0.78

constituencies, on the basis of which the regression parameters were estimated, is not a random sample, and hence might bias the forecast in the same systematic way which determines the speed of declaration. The added restriction might have the effect of correcting for that bias.

Finally, the pattern of convergence of the regression coefficient is traced throughout the iterations and depicted in Figures I and II. For all practical purposes, convergence is achieved in four or five steps. It continued to iterate to the 11th iteration only because our convergence criterion is strict.

7. A CONCLUDING REMARK

It was shown that the linear hierarchical model yields a richer family of estimates than the Brown-Zidek MRRE, but at the cost of requiring more information in order to identify all the necessary parameters in the case of unknown covariances. This results from the fact that Brown-Zidek's MRRE is an Empirical Bayes type estimate and hence all the parameters can be estimated from the sample itself. Ours, on the other hand, is a Bayes type estimate which requires some prior knowledge. Hence, the relative merits of both estimates must take into account the availability and reliability of the extra information required for the implementation of the latter.

Figure I: Convergence Pattern of B, the Matrix of Coefficients

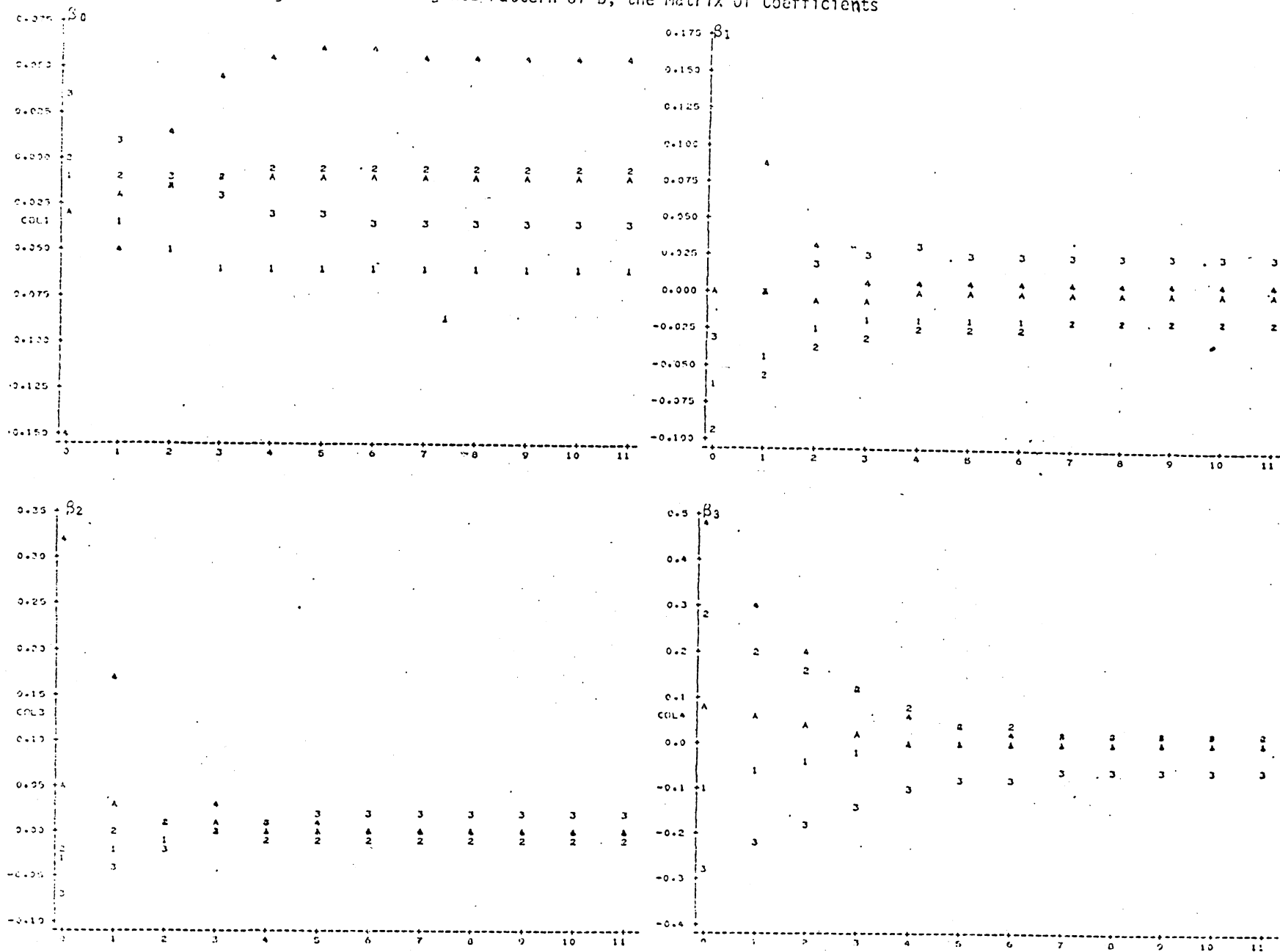
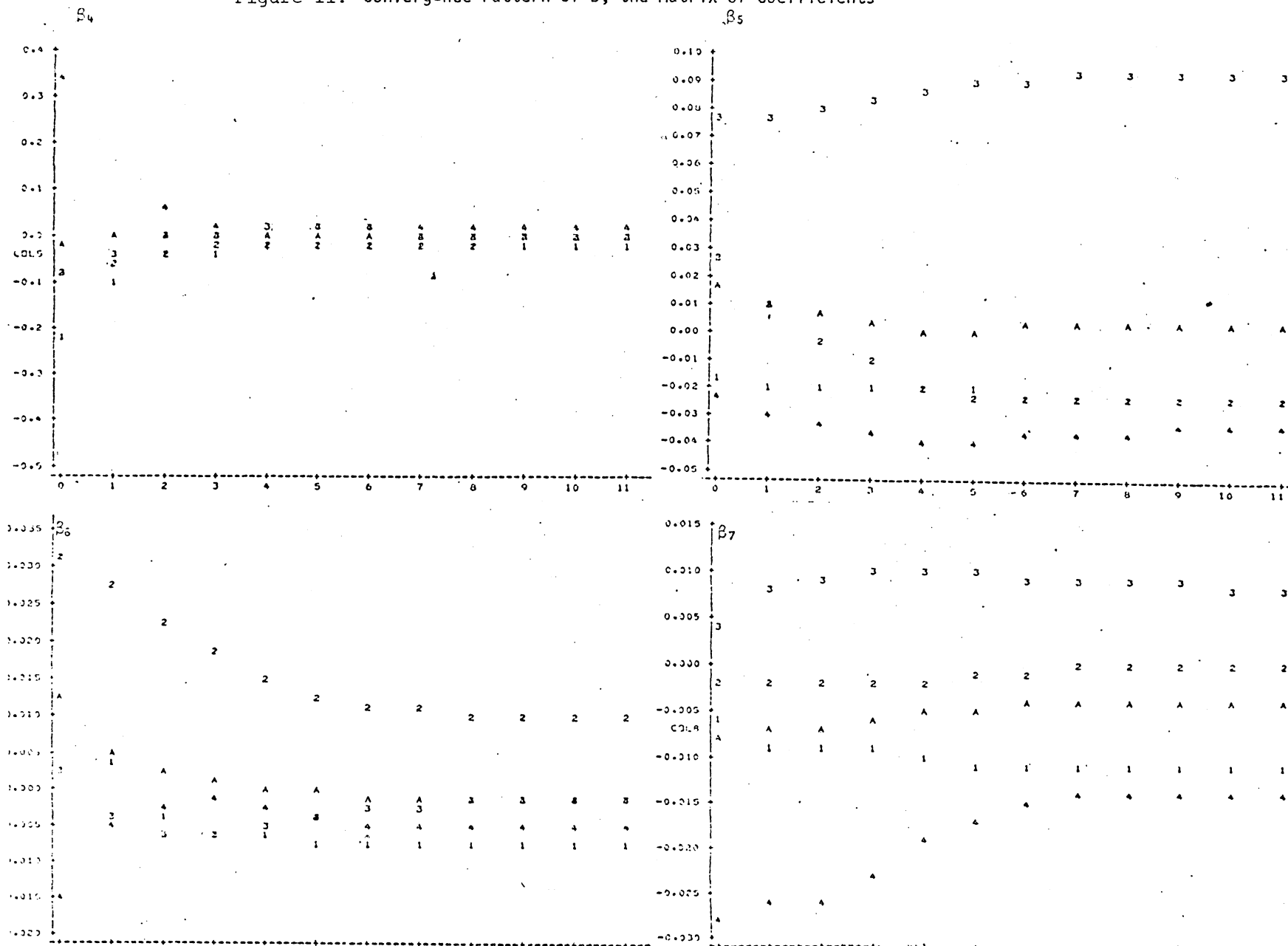


Figure II: Convergence Pattern of B, the Matrix of Coefficients



ACKNOWLEDGEMENTS

This research was supported by NSF Grant SOC 79-09674 to the University of Southern California, where the research has begun and by the ZWO to the University of Amsterdam, Faculty of Actuarial Science and Econometrics, where the research has been completed.

I am very grateful to James V. Zidek, who stimulated my interest in the problem, to Professors Charles Stein and Lawrence Brown for their very valuable comments, to Tom Wansbeek, and to many of my colleagues at the Department of Economics at USC, especially Robert Kalaba, Arie Kapteyn, and Esafandiar Maasoumi for their comments and discussions of an earlier draft.

REFERENCES

- Brown, P.J. and Payne, C. (1975). Election night forecasting (with discussion). J.R. Statist. Soc. A, 138, 463-483.
- Brown, P.J. and Zidek, J.V. (1979). Multivariate ridge regression with unknown covariance matrix. T.R. No. 79-11, Dept. of Mathematics, University of British Columbia, Vancouver, Canada.
- Brown, P.J. and Zidek, J.V. (1980). Adaptive multivariate ridge regression. Ann. Statist., 6, 64-74.
- Brown, P.J. and Zidek, J.V. (1982). Multivariate regression shrinkage estimators with unknown covariance matrix. Scand. J. Statist., 9, 209-215.
- Dempster, A.P., Rubin, D.B. and Tsutakawa, R.K. (1979). Estimation in covariance components models. J. Am. Statist. Assoc., 76, 341-353.
- Efron, Bradely and Morris, Carl (1972). Empirical Bayes on vector observations: An extension of Stein's method. Biometrika, 59, 2, 335-347.
- Fearn, T. (1975). A Bayesian approach to growth curves. Biometrika, 62, 89-100.
- Haitovsky, Y. (1979). The linear hierarchical model: I. Estimation. Modelling Research Group, Dept. of Economics, USC, MRG 7902.
- Haitovsky, Y. (1985). The linear hierarchical model and its applications in econometric analysis. In: Bayesian Inference and Decision Techniques with applications: Essays in Honor of Bruno de Finetti (ed. P.K. Goel & Arnold Zellner), forthcoming.
- Haitovsky, Y. and Wax, Y. (1980). Generalized ridge regression, least squares with stochastic prior information and Bayesian estimators. Applied Mathematics and Computation, 7, 125-154.
- Harville, D.A. (1977). Maximum likelihood approaches to variance component estimation and to related problems. J. Am. Statist. Assoc., 72, 320-338.
- Hoerl, A.E. and Kennard, R.W. (1970). Ridge regression: biased estimation for nonorthogonal problems. Technometrics, 12, 55-67.
- Lindley, D.V. and Smith, A.F.M. (1972). Bayes estimates for the linear model (with discussion). J.R. Statist. Soc., B, 34, 1-41.
- O'Hagan, A. (1976). On posterior joint and marginal modes. Biometrika, 63, 329-333.
- Sclove, S.L. (1971). Improved estimation of parameters in multivariate regression. Sankhya, A, 33, 61-66.
- Thisted, Ronald A. (1970). Ridge Regression, Minimax Estimation and Empirical Bayes Methods. PhD thesis, Stanford Un., Dept. of Statistics.
- Zellner, Arnold (1962). An efficient method of estimating seemingly unrelated regressions and test for aggregation bias. J. Am. Statist. Assoc., 57, 348-368.

111

111

100

100