

Danny Pfeiffermann, Hebrew University and Lisa LaVange, Research Triangle Institute

1. INTRODUCTION

This paper discusses the use of stochastic regression coefficients (SRC) models for analyzing survey data. These models permit different regression coefficients in different small groups of the population, an important advantage when analyzing large heterogeneous populations. The model groups can often be expected to overlap with the design clusters which by their nature are homogeneous with respect to economic and socio-demographic characteristics. The variation of the group regression coefficients can then be modelled as a function of known characteristics of the design groups. This has an important impact on the inference process and allows the problem of design selection bias to be dealt with more easily. The special features of the design require some alterations and modification to classical inference procedures but also offer some additional inference methods which employ the known selection probabilities.

The analysis of survey data using SRC models is currently being investigated at the Research Triangle Institute (RTI) under an ongoing research project with the National Institute of Child Health and Human Development (NICHD). The data base used for the analysis consists of data from the National Health and Nutrition Examination Survey II (NHANES II), conducted between 1976-1980, complemented with 1980 U.S. Census data.

The sampling design used for the NHANES II is discussed in the next section. This design is typical of many other designs used in large scale surveys and motivates some of the analysis procedures discussed in subsequent sections. Section 3 deals with the estimation of the regression coefficients and their expectations, and Section 4 deals with the estimation of the unknown variances. Empirical results of the NHANES analysis are presented in Section 5 to illustrate some of the procedures discussed in previous sections.

2. THE NHANES DESIGN

A detailed description of the NHANES is given in McDowell et al. (1981). The NHANES consists of a stratified multistage probability cluster sample of households in the U.S. At the first stage, primary sampling units (PSUs) that coincide with counties or groups of contiguous counties and selected for the National Health Interview Survey were stratified into 64 "superstrata" based on size, income and racial distribution. One PSU was selected from each stratum with probability proportional to size (PPS). In order to oversample persons with low incomes, Enumeration Districts (EDs) within the selected PSUs were sorted into poverty and non-poverty strata. EDs were then selected separately within each stratum with PPS. Households within EDs were clustered into segments of eight adjacent addresses, and a systematic sample of segments was selected across all the EDs with no more than one segment per ED. At the final stage, persons were

selected one per household, roughly, with young and old age groups being oversampled.

3. MODEL AND ESTIMATORS

In what follows we refer to the population groups with different vectors of coefficients as EDs because these are the groups used in the empirical study. Let  $Y_{ij}$  be the value of the dependent variable for unit  $j$  of ED  $i$  and  $X_{ij}$  the corresponding values of  $(k+1)$  independent variables so that  $X_{ij1} \equiv 1$ . The vector values of  $Y$  and the design matrix  $X$  observed for ED  $i$  will be denoted by  $Y_i$  and  $X_i$ . The orders of  $Y_i$  and  $X_i$  are  $(m_i \times 1)$  and  $(m_i \times (k+1))$ . It is not assumed that  $m_i > k+1$  and often  $m_i = 1$ .

3.1 The Model

$$Y_i = X_i \beta_i + \epsilon_i ; E(\epsilon_i) = 0 , E(\epsilon_i \epsilon_i^T) = \sigma^2 I \quad (3.1)$$

$$\beta_i = A_i \gamma + \eta_i ; E(\eta_i) = 0 , E(\eta_i \eta_i^T) = \Lambda \quad (3.2)$$

where  $\Lambda = \text{diag}[\delta_0^2, \delta_1^2, \dots, \delta_k^2]$ . It is assumed also that  $E(\epsilon_i \eta_i^T) = 0$  and that residuals pertaining to different EDs are independent. The matrix  $A_i$  is of order  $[(k+1) \times (k+1)(p+1)]$  and has the form  $A_i = I \otimes a_i$  where  $\otimes$  denotes the Kronecker product and  $a_i$  represents the ED characteristics with  $a_{i1} \equiv 1$ . (In the NHANES example, typical ED characteristics are income and education levels, race composition, etc. See Section 5 for the exact definitions used in the analysis).

The vector  $\gamma^T = (\gamma_0^T, \gamma_1^T, \dots, \gamma_k^T)$  is of order  $[k+1][p+1]$  so that different vectors of coefficients  $\gamma_j$  are assumed for different components  $\beta_{ij}$  of the vectors  $\beta_i$ . In a special case  $A_i = I$  implying that the coefficients  $\beta_{ij}$ ,  $i=1,2,\dots$  can be considered as random drawings from some common distribution with mean  $\gamma_j$  and variance  $\delta_j^2$ .

Inserting (3.2) into (3.1) yields the following model for the sample observations:

$$\left. \begin{aligned} Y_i &= X_i A_i \gamma + X_i \eta_i + \epsilon_i = X_i^* \gamma + \mu_i \\ E(\mu_i) &= 0; E(\mu_i \mu_i^T) = \sigma^2 I + X_i \Lambda X_i^T; E(\mu_i \mu_k^T) = 0 \quad i \neq k \end{aligned} \right\} \quad (3.3)$$

Notice that the  $l$ -th row of the matrix

$$\left. \begin{aligned} X_i^* &= [1, \dots, 1] \text{ is} \\ X_{i1}^{*T} &= [a_i^T, x_{i11}^T a_i^T, \dots, x_{i1k}^T a_i^T] \end{aligned} \right\} \quad (3.4)$$

Thus, the model defined by (3.3) can be considered as an extension of the model defined by (3.1) obtained by adding the regressors included in A and the first-order (multiplicative) interactions between the regressors in  $X_i$  and A. Furthermore, in many applications the regressors in A are of similar nature to those in  $X_i$  so that to some extent (3.3) represents a polynomial regression equation. The model for a sample of n EDs can be written compactly as

$$\underline{y} = X^* \underline{\gamma} + \underline{\mu}; E(\underline{\mu}) = \underline{0}, E(\underline{\mu} \underline{\mu}^T) = V \quad (3.5)$$

where  $\underline{y}^T = (y_1^T \dots y_n^T)$ ,  $X^{*T} = [X_1^{*T} \dots X_n^{*T}]$ ,

$\underline{\mu}^T = (\mu_1^T \dots \mu_n^T)$  and V is block diagonal with

$$V_{ii} = E(\mu_i \mu_i^T).$$

### 3.2 Estimation of $\underline{\gamma}$

The matrix  $X^*$  would usually be of full rank  $P^* = (P+1)(k+1)$ . For known variances  $\sigma^2$  and  $\delta_j^2$ , the best linear unbiased (BLUE) generalized least squares (GLS) estimator of  $\underline{\gamma}$  under the model is

$$\hat{\underline{\gamma}}_{GLS} = (X^{*T} V^{-1} X^*)^{-1} X^{*T} V^{-1} \underline{y} \quad (3.6)$$

In practical applications the variances are unknown and have to be estimated from the sample. Inserting these sample estimators into (3.6) yields the empirical estimators  $\hat{\underline{\gamma}}_E$ .

Assuming that the variance estimators are consistent, the asymptotic distribution of  $\hat{\underline{\gamma}}_E$  is under general conditions the same as that of  $\hat{\underline{\gamma}}_{GLS}$ , Anderson (1973).

The OLS estimator of  $\underline{\gamma}$  is

$$\hat{\underline{\gamma}}_{OLS} = (X^{*T} X^*)^{-1} X^{*T} \underline{y} \quad (3.7)$$

Under the model  $\hat{\underline{\gamma}}_{OLS}$  is unbiased and under mild conditions it is consistent.

The performance of  $\hat{\underline{\gamma}}_{BLUE}$  and  $\hat{\underline{\gamma}}_{OLS}$  depends on the noninformativeness of the design. A design is noninformative if the selection of the sample is independent of the model regression residuals. Under an informative design, the model holding for the sample data is different from the model holding in the population, and in such cases estimators like  $\hat{\underline{\gamma}}_{GLS}$  and  $\hat{\underline{\gamma}}_{OLS}$  can be severely biased. This may happen in the NHANES case, for example, if not all of the design characteristics determining the selection of EDs are included in the vectors  $\underline{a}_i$ , so that

$f(\underline{\beta}_i | \underline{a}_i, i \in S_E) \neq f(\underline{\beta}_i | \underline{a}_i)$  where  $S_E$  denotes the sample of EDs. See e.g., Holt Smith and Winter (1980) for discussion and references on the notion of informative designs.

A simple way to deal with the informativeness of the design is to weight every observation

vector  $(y_{ij}, x_{ij})$  by the inverse of its selection probability. The weighted least squares (WLS) estimator can be written as

$$\hat{\underline{\gamma}}_{WLS} = (X^{*T} W X^*)^{-1} X^{*T} W \underline{y} \quad (3.8)$$

where W is a diagonal matrix with the sampling weights on the main diagonal. The estimator  $\hat{\underline{\gamma}}_{WLS}$  is again unbiased and consistent under the model. However, its main property is that it is approximately design ("P") unbiased and consistent for the Census vector  $\underline{\gamma}_{CEN}$  which in turn is model ("ξ") unbiased and consistent for  $\underline{\gamma}$ . The census vector is the estimator that would have been obtained in case of a census. For a definition of consistency in finite population sampling, see e.g., Isaki and Fuller (1982).

It follows from the discussion above that  $\hat{\underline{\gamma}}_{WLS}$  is approximately unbiased and consistent for  $\underline{\gamma}$  with respect to the  $P\xi$  distribution. The variance of  $\hat{\underline{\gamma}}_{WLS}$  can be decomposed as

$$\begin{aligned} \text{Var}_{P\xi}(\hat{\underline{\gamma}}_{WLS}) &= E_{\xi}[\text{Var}_P(\hat{\underline{\gamma}}_{WLS})] + \\ &+ \text{Var}_{\xi}[E_P(\hat{\underline{\gamma}}_{WLS})] = E_{\xi}[\text{Var}_P(\hat{\underline{\gamma}}_{WLS})] + 0[M^{-1}] \end{aligned} \quad (3.9)$$

where M denotes the population size. Eq. (3.9) implies that in practical situations where the sampling fractions  $m/M$  are very small, estimators of  $\text{Var}_{P\xi}(\hat{\underline{\gamma}}_{WLS})$  can be obtained by estimat-

ing  $\text{Var}_P(\hat{\underline{\gamma}}_{WLS})$ . Such estimators are calculated for general multistage probability sampling designs by the commonly used software packages for regression analysis of survey data such as SURREGR (Holt 1977), OSIRIS (1979) and SUPERCARP (Hidirolou et. al, 1980).

### 3.3 Prediction of $\underline{\beta}_i$ When Variances are Known

The use of SRC models requires the specification of the target vectors of coefficients. In general, these will be of the form

$$\underline{\beta}_C = \sum_{i=1}^N C_i \underline{\beta}_i \text{ with } \sum_{i=1}^N C_i = 1 \text{ where "N" denotes}$$

the number of population EDs. Simple examples are (i)  $C_i = 1/N$  (ii)  $C_i = M_i/M$  where  $M_i$  is the size of ED i and (iii)  $C_i = \frac{1}{T}$  for  $i \in U_T$ ,  $C_i = 0$  otherwise, where  $U_T$  is a subgroup (domain) of T EDs defined by the ED characteristics. For

known variances  $\sigma^2$  and  $\delta_j^2$ ,  $j=0..k$ , the BLUE predictor of  $\underline{\beta}_C$  under the model is

$$\begin{aligned} \hat{\underline{\beta}}_C &= \sum_{i=1}^N C_i \hat{\underline{\beta}}_i \text{ where for } i \in S_E, \\ \hat{\underline{\beta}}_i &= A_i \hat{\underline{\gamma}}_{GLS} + \Delta X_i^T V_{ii}^{-1} (y_i - X_i^* \hat{\underline{\gamma}}_{GLS}) \end{aligned} \quad (3.10)$$

and  $\hat{\beta}_i = A_i \hat{\gamma}_{GLS}$  otherwise (Pfeffermann, 1984).

The distributional properties of  $\hat{\beta}_i$  are with respect to the joint  $\xi$ -distribution of the  $\beta_i$ 's and the  $Y_{ij}$ 's given the selected sample.

Notice from (3.10) that the optimal estimator of  $\beta_i$  for  $i \in S_E$  is the same as the optimal estimator of its expectation. For the sampled EDs these estimators are corrected by taking into account the deviations of the observations  $Y_{ij}$  for their expected means. In cases where the target vector  $\beta_C$  is a population average, these correction factors would usually have a minor effect on the estimator. Thus, an alternative predictor of  $\beta_C$  in such situations is

$\tilde{\beta}_C = \sum_{i=1}^N C_i A_i \hat{\gamma}$  where  $\hat{\gamma}$  can be any one of the estimators discussed in the previous section.

The estimators considered so far assume that the weighted averages  $C_A = \sum_{i=1}^N C_i A_i$  are known.

When this is not the case, they can be estimated using the sampling weights, i.e.,

$$\hat{C}_A = \sum_{i \in S_E} C_i A_i W_i / \sum_{i \in S_E} W_i. \quad (3.11)$$

The estimator  $\hat{C}_A$  is an example of the use of the design in an essentially model based analysis.

#### 4. ESTIMATION OF THE UNKNOWN VARIANCES

##### 4.1 Maximum Likelihood Estimators

Estimation of the unknown variances  $\sigma^2$  and  $\delta_j^2$ ,  $j=0, \dots, k$  is needed for three main reasons: (i) for calculating the empirical estimator  $\hat{\gamma}_E$  and the correction factors defined by equation (3.10), (ii) for testing hypotheses regarding the variation of the vectors of coefficients, and (iii) for constructing confidence intervals for the unknown coefficients. Variance estimation in SRC models, or more generally, in variance component models has been a major area of research in recent years. Methods with known properties and feasible computational algorithms have been developed and they are widely discussed in the literature (see e.g., Searle, 1971, Harville, 1977, and the more recent article by Henderson, 1984).

In this section we discuss possible ways of deriving maximum likelihood estimators (m.l.e.) for the unknown variances. In particular we emphasize that m.l.e. can be obtained by repeated use of regression software routines without the need for further software development. This has some additional advantages which are outlined following the description of the computations.

The results presented below are borrowed from results on variance estimation in the 'classical' variance components model for which the V-C matrix of the residuals is linear in the unknown parameters. The general model defined

by (3.1) and (3.2) falls under that category as can be seen by rewriting equations (3.3) in the form

$$Y = X^* \gamma + U_0 \xi_0 + U_1 \xi_1 + \dots + U_k \xi_k + \xi \quad (4.1)$$

where  $U_j$  is a block diagonal matrix of order  $(m \times n)$  with the  $i$ -th block being the  $j$ -th column  $X_{i(j)}$  of the design matrix  $X_i$ , and  $\xi_j = (\eta_{1j}, \eta_{2j}, \dots, \eta_{nj})$ . The vector  $\xi_j$  consists of the deviations  $(\beta_{1j} - \bar{\beta}_1^T \gamma_j)$  of the regression coefficients in the various groups so that  $\text{Var}(\xi_j) = \delta_j^2 I_n$  and  $E(\xi_i \xi_j^T) = 0$  for  $i \neq j$ . It follows that

$$\text{Var}(Y) = \sum_{j=0}^k \delta_j^2 U_j U_j^T + \sigma^2 I_m = V \quad (4.2)$$

The model defined by (4.1) and (4.2) is the classical mixed model of variance components. Assuming that the error terms  $\xi_j$  and  $\xi$  have a normal distribution, the likelihood equations are (Anderson, 1973),

$$\gamma = (X^{*T} V^{-1} X^*)^{-1} X^{*T} V^{-1} Y \quad (4.3)$$

$$B \hat{\xi}^* = \xi \quad (4.4)$$

where

$$B = (B_{pq}), \quad B_{pq} = \text{trace} [V^{-1} (U_p U_p^T) V^{-1} (U_q U_q^T)]$$

$$\xi = (c_p), \quad (c_p) = [(Y - X^* \gamma)^T V^{-1} U_p U_p^T V^{-1} (Y - X^* \gamma)]$$

$$\text{and } \hat{\xi}^{*T} = (\delta_0^T \dots \delta_k^T, \sigma^2).$$

The common statistical software packages such as SAS and BMDP do not include procedures for solving the likelihood equations (4.4) for the general case where the elements of the matrices  $X^*$  and  $U_j$  are different from zeroes and ones. A SAS procedure which puts out a SAS data set containing the coefficients for the system (4.4) for given (prior) values of the elements of  $V$  has been written at North Carolina State University (Giesbrecht, 1985). Iterating on this system of equations with newly obtained estimators for  $\hat{\xi}^*$  used as prior values for defining  $V$  in the next iteration yields m.l.e. of the unknown variances provided the system converges and no negative values are encountered.

An alternative procedure to obtain m.l.e. is to transfer the problem into a problem of generalized regression analysis. This is done by fitting a linear model to the squares and cross products of the residuals

$$\mu_{ij} = (Y_{ij} - X_{ij}^* \gamma)^2, \text{ i.e.,}$$

$$\begin{aligned}
E(\mu_{ij}^2) &= \delta_0^2 + x_{ij1}^2 \delta_1^2 + \dots + x_{ijk}^2 \delta_k^2 + \sigma^2 = \\
&= (\tilde{z}_{ij}^T, 1) \tilde{\xi}^* \\
E(\mu_{ij} \mu_{il}) &= \delta_0^2 + x_{ij1} x_{il1} \delta_1^2 + \dots \\
&\quad + (x_{ijk} x_{ilk}) \delta_k^2 = (\tilde{z}_{ijl}^T, 0) \tilde{\xi}^* \quad j \neq l \\
E(\mu_{ij} \mu_{i^*j^*}) &= 0 \quad \text{for } i \neq i^*
\end{aligned} \tag{4.5}$$

$$\begin{aligned}
\text{Var}(\mu_{ij}^2) &= 2(\tilde{x}_{ij}^T \Lambda \tilde{x}_{ij})^2 \\
\text{Cov}(\mu_{ij} \mu_{il}, \mu_{ip} \mu_{iq}) &= \begin{pmatrix} \tilde{x}_{ij}^T \Lambda \tilde{x}_{ip} \\ \tilde{x}_{il}^T \Lambda \tilde{x}_{iq} \end{pmatrix} \begin{pmatrix} \tilde{x}_{il}^T \Lambda \tilde{x}_{iq} \\ \tilde{x}_{ij}^T \Lambda \tilde{x}_{ip} \end{pmatrix} \\
&\quad + \begin{pmatrix} \tilde{x}_{ij}^T \Lambda \tilde{x}_{iq} \\ \tilde{x}_{il}^T \Lambda \tilde{x}_{ip} \end{pmatrix} \begin{pmatrix} \tilde{x}_{il}^T \Lambda \tilde{x}_{iq} \\ \tilde{x}_{ij}^T \Lambda \tilde{x}_{ip} \end{pmatrix} \\
\text{Cov}(\mu_{ij} \mu_{il}, \mu_{i^*j^*} \mu_{i^*l^*}) &= 0 \quad \text{for } i \neq i^*
\end{aligned} \tag{4.6}$$

Equations (4.5) follow directly from (3.3). Equations (4.6) translate equations (19) of Anderson (1973) to the present model. Let  $\tilde{\mu}_i$  be the vector of squares and cross products of the residuals  $\mu_{ij}$  corresponding to ED  $i$  and arranged in some convenient order, and define  $\tilde{\mu}^T = (\tilde{\mu}_1^T \dots \tilde{\mu}_n^T)$ . Let  $Z_i$  be the matrix consisting of the rows  $(\tilde{z}_{ij}^T, 1)$  and  $(\tilde{z}_{ijl}^T, 0)$  arranged in the same order as the elements of  $\tilde{\mu}_i$  and define  $Z^T = (Z_1^T \dots Z_n^T)$ . It follows from (4.5) and (4.6) that  $\tilde{\mu}_i = Z_i \tilde{\xi}^* + \tilde{\epsilon}_i$ ,  $E(\tilde{\epsilon}_i) = 0$ ,  $E(\tilde{\epsilon}_i \tilde{\epsilon}_i^T) = \Lambda_i$  where  $\Lambda_i$  is defined by (4.6). The linear model holding for  $\tilde{\mu}$  is

$$\tilde{\mu} = Z \tilde{\xi}^* + \tilde{\epsilon}; \quad E(\tilde{\epsilon}) = 0, \quad E(\tilde{\epsilon} \tilde{\epsilon}^T) = \Lambda = \text{diag}[\Lambda_1 \dots \Lambda_n] \tag{4.7}$$

It follows from (4.7) that for a given, 'known' V-C matrix  $\Lambda$ , the BLUE estimators of the unknown variances are obtained as

$$\tilde{\xi}^* = (Z^T \Lambda^{-1} Z)^{-1} Z^T \Lambda^{-1} \tilde{\mu} \tag{4.8}$$

Notice from (4.5) that as long as some of the ED's contain more than one observation,  $\delta_0^2$  and  $\sigma^2$  can be estimated separately.

Anderson (1973) has shown that m.l.e. of  $\tilde{\xi}^*$  can be obtained by iterating between (4.3) and (4.8) with  $\tilde{\mu}$  replaced by  $(\tilde{y} - X^* \hat{\gamma})$ . (See also Brown and Burgess, 1984 for further discussion and examples.) Thus, an alternative to solving the system of equations in (4.4) is to solve equation (4.8) with  $\tilde{\mu} = (\tilde{y} - X^* \hat{\gamma})$ .

It is important to emphasize that the cross-products of residuals  $\mu_{ij}$  and  $\mu_{i^*j^*}$  pertaining to different ED's  $i$  and  $i^*$  are not considered in the model (4.7). This is so because for  $i \neq i^*$

$$E(\mu_{ij} \mu_{i^*j^*}) = \text{Cov}(\mu_{ij} \mu_{i^*j^*}, \mu_{pq} \mu_{pq^*}) = 0 \tag{4.9}$$

Thus, the set of cross-products  $\{\mu_{ij} \mu_{i^*j^*}, i \neq i^*\}$  have zero expectations and they are uncorrelated with the squares and cross-products of residuals included in the vector  $\tilde{\mu}$ . Their inclusion in the model (4.7) has therefore no impact on the estimation of  $\tilde{\xi}^*$ .

The immediate implication of (4.9) is that for situations in which there are many ED's with only a few observations in each, employing the model (4.7) and iterating between (4.3) and (4.8) is a simple operation. Notice in this respect that

$$\begin{aligned}
V^{-1} &= \text{diag}[V_1^{-1} \dots V_n^{-1}], \quad V_i^{-1} = (\sigma^2 I_{m_i} + X_i \Lambda X_i^T)^{-1} \\
&= \frac{1}{\sigma^2} [I - X_i (X_i^T X_i + \sigma^2 \Lambda^{-1})^{-1} X_i^T]
\end{aligned} \tag{4.10}$$

so that the computation of  $\hat{\gamma}$  and  $\hat{\xi}^*$  involves only the inversion of matrices of orders  $(k+1)$ ,  $(k+2)$  and  $(k+1)(p+1)$ .

#### 4.2 Estimation of Variances of Variance Estimators Under Informative and Noninformative Designs

Estimators for the model based variances of the variance estimators are readily obtained from (4.8) and are given by

$$\hat{\text{Var}}(\hat{\xi}^*) = (Z^T \hat{\Lambda}^{-1} Z)^{-1} \tag{4.11}$$

where  $\hat{\Lambda}$  denotes the estimator of  $\Lambda$  obtained in the last iteration of the procedure described above.

Design based estimators of the variances and covariances of  $\hat{\xi}^*$  can be obtained using a survey regression software package such as SURREGR (Holt, 1977). The use of such estimators becomes essential when the design is informative, in which case  $\hat{\xi}^*$  has to be replaced by the probability weighted estimator in each stage of the analysis.

Another important advantage of the regression method for deriving m.l.e. is that it permits a simple way of constraining the estimates to nonnegative values. An unconstrained solution of the likelihood equations can yield negative variance estimators. Instead of solving (4.8), one can solve the minimization problem

$$\min_{\tilde{\xi}^* \geq 0} \left[ (\tilde{\mu} - Z \tilde{\xi}^*)^T \Lambda^{-1} (\tilde{\mu} - Z \tilde{\xi}^*) \right] \tag{4.12}$$

By iterating between (4.3) and (4.12) with the most recently obtained estimators of  $\hat{\xi}^*$  inserted

in V and A to obtain new estimators, the solution of (4.12) reduces to a problem of minimizing a quadratic function subject to linear inequality constraints for which there are computer programs available. Brown and Burgess (1984) discuss the use of constrained variance estimators and give empirical results which illustrate important features of the approach.

## 5. EMPIRICAL STUDY

### 5.1 Statement of Problem

The method of fitting stochastic regression coefficients models was applied to NHANES II data to assess the relationship between blood lead and blood pressure. Several articles appearing recently in the literature have dealt with this problem using data from both the NHANES I and NHANES II surveys. Pirkle, et al. (1985) found a significant relationship between blood lead and both systolic and diastolic blood pressure in white males aged 40-59 years while controlling for other dietary and anthropometric blood pressure correlates. Harlan, et al. (1985) also found blood lead to be a significant predictor of both systolic and diastolic blood pressure in men aged 12 to 74 years, but not in women. Shaper and Pocock (1985) used data from the British Regional Heart Study and in contrast to the other studies found no significant relationship between these two variables in middle-aged men.

Because the main goal of the empirical study was to investigate methodological issues as opposed to deriving the "best" explanatory model, it was decided to begin with previously published model equations, thereby taking advantage of the extensive variable selection procedures used to derive these models. Analysis was thus carried out on the subgroup consisting of all males aged 12-74 years using Harlan's model as a basis.

### 5.2 Model Definitions and Results

In order to fit SRC models to NHANES data, it was decided to define population subgroups according to EDs, (the design clusters selected at the second stage of the NHANES sample) for two reasons: (i) EDs are generally thought to be homogeneous clusters of households, and (ii) design information and selection probabilities, were available for each ED. This information had been merged from the 1970 U.S. Census tapes onto the NHANES II sample of household segments.

As explained above, the model for diastolic blood pressure derived by Harlan, et al. (1985) for males aged 12-74 years was used as a basis for fitting stochastic regression coefficients models to the NHANES data. This model was refit to the sample after deleting observations for which the ED affiliation was missing (n = 2676 versus n = 2818 in Harlan's analysis). The model results are presented in column 5 of Table 1. The coefficients were estimated by probability weighted least squares using RTI's survey regression software package SURREGR (Holt, 1977).

SRC models as defined by equations (3.1) and (3.2) were fit to the data. The vectors of ED specific regression coefficients were modelled as functions of the following ED characteristics: medium family income,

proportion of non-whites, proportion of persons aged 25+ with less than 9 years of education, urban/rural living areas, and a poverty stratum indicator. The domain of analysis consisted of 1,892 EDs.

The significant  $\gamma$  coefficients of the models identified for the regression coefficients (equation 3.2) are presented in columns 2-4 of Table 1. As can be seen, a significant portion of the variation among the ED specific regression coefficients for age, age<sup>2</sup>, race and particularly blood lead is explained by two ED characteristics; the proportion of non-whites and the proportion of persons with education of less than 9 years of school.

The last column of the table was calculated as

$$\hat{\beta}_C = \left( \sum_{i=1}^n M_i \hat{\beta}_i / \Pi_i \right) / \left( \sum_{i=1}^n M_i / \Pi_i \right) \text{ where } \{M_i\}$$

and  $\{\Pi_i\}$  are the ED sizes and selection probabilities respectively. The entries of this column estimate the average regression relationship in the finite population and are found to be very similar to the coefficients obtained by Harlan's model presented in Table 1.

### Summary

The results obtained at this stage suggest that whereas a unique regression equation can be assumed to assess the average effect of the regressors included in the original model, such a model may be inappropriate for studying the relationships in subdomains of the population. Obviously, a more extensive analysis is needed to support these early suggestions. We have started programming the estimation of the unknown variances using the procedures described in Section 4. At the same time, we are trying to locate some additional variables to be incorporated in the regression equation explaining the variation of the regression coefficients among the population EDs.

### REFERENCES

- Anderson, T.W. (1973). "Asymptotically Efficient Estimation of Covariance Matrices with Linear Structure." *The Annals of Statistics*, Vol. 1, 135-141.
- Brown, K.G. and Burgess, M.A. (1984). "On Maximum Likelihood and Restricted Maximum Likelihood Approaches to Estimation of Variance Components." *J. Statist. Comput. Simul.*, Vol. 4, 1-18.
- Giesbrecht, F. (1985). "MIXMOD, A SAS Procedure for Analyzing Mixed Models. Mimeograph Series No. 1659, North Carolina State University, Raleigh, NC.
- Harlan, W.R., Landis, J.R., Schmouder, R.L., Goldstein, N.G. and Harlan, L.C. (1985). "Blood Lead and Blood Pressure, *JAMA*, 253, 530-534.
- Harville, D.A. (1977). "Maximum Likelihood Approaches to Variance Component Estimation and to Related Problems." *J. Amer. Statist. Ass.*, 72, 320-340.

- Henderson, C.R. (1984). "ANOVA, MINQUE, REML and ML Algorithms for Estimation of Variances and Covariances." Proceedings of the 50th Anniversary Conference, Iowa State Statistical Laboratory, H. David and H.T. David, Editors.
- Hidiroglou, M., Fuller, W.A. and Hickman, R.D. (1980). "SUPER CARP." Statistical Laboratory, Iowa State University, Ames, Iowa.
- Holt, D., Smith, T.M.F. and Winter, P.D. (1980). "Regression Analysis of Data from Complex Surveys." *J.R. Statist. Soc. A*, 143, 474-487.
- Holt, M.M. (1977). SURREGR: Standard Errors of Regression Coefficients from Sample Survey Data." Research Triangle Institute, Research Triangle Park, NC.
- Isaki, C.T. and Fuller, W.A. (1982). "Survey Design Under the Regression Superpopulation Model." *J. Amer. Statist. Ass.*, 77, 89-96.
- McDowell, A., Engel, A., Massey, J.T. and Maurer, K. (1981). "Plan and Operation of the Second National Health and Nutrition Examination Survey 1976-80." DHHS Publication No. (PHS) 81-1317, U.S. Department of Health and Human Services, National Center for Health Statistics, Hyattsville, MD.
- OSIRIS (1984). Survey Research Center Computer Support Group, The Institute for Social Research, The University of Michigan, Ann Arbor, MI.
- Pfeffermann, D. (1984). "On Extension of the Gauss-Markov Theorem to the Case of Stochastic Regression Coefficients." *J.R. Statist. Soc. B*, 46, 139-148.
- Pirkle, J.L., Schwartz, J., Landis, J.R. and Harlan, W.R. (1985). "The Relationship Between Blood Lead Levels and Blood Pressure and Its Cardiovascular Risk Implications." *American Journal of Epidemiology*, 121, 246-258.
- Searle, S.R. (1971). "Topics in Variance Component Estimation." *Biometrics* 27, 1-76.
- Shaper, A.G. and Pocock, S.J. (1985). "Blood Lead and Blood Pressure." *British Medical Journal*, 291, 1147-1149.

TABLE 1. MODELS IDENTIFIED FOR DISTINCT COEFFICIENTS, ESTIMATES OF REGRESS COEFFICIENTS

$\beta_j$ Coeff.	Models Identified <sup>1</sup>			Est. of Coefficients	
	Constant	Prop. Non. W	Prop. Edu. < 9	Unique Reg Line <sup>2</sup>	Weighted Avg. of ED Coeff. <sup>3</sup>
Coefficients					
Intercept	43.486		-17.871	37.407	38.492
Age	0.449	0.987		0.524	0.539
Age <sup>2</sup>	- 0.004	-0.010		- 0.005	- 0.005
Body mass	0.960			0.949	0.960
Race			8.332	- 0.812	2.328
Race x Age	-	-	-	0.091	-
Lead (ln)		-7.241	6.298	1.378	1.095
Hemoglobin	1.207			1.258	1.207
Serum zinc (ln)	- 3.263			- 3.223	- 3.263

<sup>1</sup>Table shows significant coefficients at the 0.07 level in the equation

$$E(\beta_{ik}) = \gamma_{0k} + a_{i1} \gamma_{1k} + \dots + a_{ip} \gamma_{pk}$$

<sup>2</sup>Model published by Harlan et al. (1985).

$$\text{<sup>3</sup>Coefficients estimated as } \hat{\beta}_i = \frac{\sum_1^n M_i A_i \hat{\gamma}_{WLS}/\Pi_i}{\sum_1^n M_i/\Pi_i}$$