Cathy Campbell, University of Minnesota

1. Introduction

In this paper we present a simple regression model for two-stage cluster samples. It is hypothesized that the model may be appropriate for many situations in which both the clusters and elements within the clusters are assumed to be sampled from infinite populations. As such, it is not a model for sampling from finite populations, but may also be considered as a superpopulation model for two-stage samples taken from finite populations.

The model of interest can be given as

 $y_{ij} = \mu + \underline{x}'_{ij}\beta + u_{ij}$ (i=1,...,b; j=1,...,n_i), or in matrix notation as

$$y = \underline{1}\mu + \underline{X}\underline{\beta} + \underline{u} \tag{1.1}$$

where <u>1</u> is an nxl vector of l's,

 $\underline{y} = (y_{11}, \dots, y_{bn_b})'$ is an nxl vector of

observed dependent variables,

- X is an nxp matrix of observed predictor variables,
- μ and β (px1) are unknown parameters
- <u>u</u> is an nxl vector of unobserved random variables with

$$E(\underline{u}|X) = 0,$$

$$Var(\underline{u}|X) = \sigma^{2}V,$$

$$V = (1-\rho_{y})I + \rho_{y}\begin{bmatrix}J_{n_{1}} & \emptyset\\ & \\ & \\ & & \\$$

$$= (1-\rho_y)I + \rho_y X_b X_b', \qquad (1.3)$$

$$J_n \text{ is an } n_i x n_i \text{ matrix of 1's,}$$

$$I_b = \begin{bmatrix} 1 & \emptyset \\ 0 & 1 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \text{ is the matrix of indica-tor variables identifying the cluster from which each element was sampled,}$$

 $\stackrel{\rho}{y}$ is the intraclass correlation of the residuals around the regression line,

$$\frac{-1}{n_b-1} < \rho_y \le 1.$$

From (1.2) and (1.3), it is clear that

$$Var(y_{ij}) = \sigma^{2}$$

$$Cov(y_{ij}, y_{ik}) = \rho_{y}\sigma^{2} \quad (j \neq k)$$

$$Cov(y_{ij}, y_{k\ell}) = 0 \quad (i \neq k)$$

Since the constant σ^2 appears as a constant multiplier on all variance expressions, and will cancel from all ratios, for convenience we assume $\sigma^2=1$.

Our interest in this model lies in studying the estimation of $\underline{\beta}$; μ is considered a nuisance parameter. Conditional on X, the weighted least squares (WLS) estimator of $\underline{\beta}$ is BLU, but is rarely used because it depends on unknown parameters and is difficult to compute. More often, because of availability of computer programs and the familiarity of the technique, ordinary least squares (OLS) is used to estimate $\underline{\beta}$. In this paper we wish to consider two aspects of the estimation of $\underline{\beta}$.

(1) Sometimes cluster samples are taken for convenience or economy, sometimes from necessity. What would be the effect on the variance of the parameter estimates if a simple random sampling procedure were used instead? In sampling terminology we wish to study the design effect for the OLS estimator of $\underline{\beta}$.

(2) Is OLS an efficient estimation procedure when model (1.1) holds? If OLS is extremely inefficient, then perhaps some form of approximate WLS, using an estimate of ρ_y , should be considered as an alternative.

For convenience, we restrict our results here to models with one or two predictor variables and consider the issue of design effects first.

2. Design Effect for Simple Linear Regression

When p=1, model (1.1) becomes

$$\underline{\mathbf{y}} = \underline{\mathbf{l}}\boldsymbol{\mu} + \underline{\mathbf{x}}\boldsymbol{\beta} + \underline{\mathbf{u}} \,. \tag{2.1}$$

We assume that \underline{x} has been transformed so that $\underline{x}'\underline{1} = 0$. Then the OLS estimator of β is given by

$$\hat{\beta}_{o} = (\underline{x}'\underline{x})^{-1}\underline{x}'\underline{y} , \qquad (2.2)$$

and

$$\operatorname{Var}(\hat{\boldsymbol{\beta}}_{o}|\underline{x}) = (\underline{x}'\underline{x})^{-2}\underline{x}' \nabla \underline{x} . \qquad (2.3)$$

Following Frankel (1971), we define the design effect of $\hat{\beta}_0$, Deff($\hat{\beta}_0$), as the ratio of the variance of $\hat{\beta}_0$ under model (1.1) to the variance of $\hat{\beta}_0$ under the assumption of a simple random selection of elements of the same overall sample size. As Var(y_{ij})= $\sigma^2 = 1$, Var($\hat{\beta}_0 | \underline{x}$) with simple random sampling is $(\underline{x}' \underline{x})^{-1}$. Therefore

$$\operatorname{Deff}(\hat{\beta}_{0}|\underline{x}) = \frac{(\underline{x}'\underline{x})^{-2}\underline{x}'\underline{v}\underline{x}}{(\underline{x}'\underline{x})^{-1}} = (\underline{x}'\underline{x})^{-1}\underline{x}'\underline{v}\underline{x} . \quad (2.4)$$

More correctly, the expression (2.4) should be called a conditional design effect since the same \underline{x} is used in both numerator and denominator.

Without loss of generality, we may assume $\underline{x'x} = 1$ in (2.4) and obtain

Deff
$$(\hat{\beta}_0 | \underline{x}) = \underline{x}' V \underline{x}$$
 $(\underline{x}' \underline{x} = 1, \underline{x}' \underline{1} = 0).$ (2.5)
Substituting (1.3) for V in (2.5) yields

$$Deff(\hat{\beta}_{o}|\underline{x}) = 1 + (\underline{x}' X_{b} X_{b}' \underline{x} - 1) \rho_{y} . \qquad (2.6)$$

To make (2.6) more easily comparable with the usual expressions for design effects, it is convenient to express $\underline{x}' X_b X_b \underline{x}'$ in terms of the intra-

class correlation, ρ_X , of the observed <u>x</u>. We first note that $\underline{x}'X_bX_b\underline{x}'$ is the sum of squares of the b cluster totals of <u>x</u>, and can be expressed as b 2 2

$$\underline{\mathbf{x}}' \mathbf{X}_{\mathbf{b}} \mathbf{X}_{\mathbf{b}}' \underline{\mathbf{x}} = \sum_{i=1}^{\infty} n_{i}^{2} \mathbf{x}_{i}^{-2}.$$
 (2.7)

Following Murthy (1967), we use the following definition of ρ_X that is applicable for unequal cluster sizes:

$$\rho_{\mathbf{X}} = \frac{\sum_{i=1}^{\mathbf{b}} \sum_{k\neq j}^{\mathbf{n}_{i}} \sum_{\mathbf{i}_{j}=\mathbf{x}}^{\mathbf{n}_{i}} (\mathbf{x}_{ij}-\mathbf{x}) (\mathbf{x}_{ik}-\mathbf{x})}{\sum_{i=1}^{\mathbf{b}} n_{i} (n_{i}-1)\sigma_{\mathbf{X}}^{2}} .$$
 (2.8)

Using the relationships $\sigma_X^2 = \frac{1}{n}$ and $\overline{x} = 0$, (2.8) reduces to

$$\rho_{\mathbf{X}} = \frac{\sum \mathbf{n}_{\mathbf{i}}^2 \overline{\mathbf{x}}_{\mathbf{i}}^2 - 1}{\frac{\sum \mathbf{n}_{\mathbf{i}}^2}{\sum \mathbf{n}_{\mathbf{i}}} - 1}$$

which gives

$$\Sigma \mathbf{n}_{i}^{2} \overline{\mathbf{x}}_{i}^{2} = 1 + \left(\frac{\Sigma \mathbf{n}_{i}^{2}}{\Sigma \mathbf{n}_{i}} - 1\right) \rho_{X} \quad (2.9)$$

9

Substituting (2.9) in (2.6) gives

$$Deff(\hat{\beta}_{0}|\underline{x}) = 1 + \left(\frac{\Sigma n_{i}^{2}}{\Sigma n_{i}} - 1\right) \rho_{X} \rho_{Y} \qquad (2.10)$$
$$= 1 + \left(\frac{Var(n_{i})}{\overline{n}} + \overline{n} - 1\right) \rho_{X} \rho_{Y} (2.11)$$

where $n = \sum_{i=1}^{\infty} n_i / b$ is the average sample size.

Noting that (2.9) is the design effect (see (2.13)) for estimating the mean of <u>x</u>, we can also obtain

$$Deff(\hat{\beta}_{o}|\underline{x}) = 1 + (Deff(\overline{x}) - 1)\rho_{Y} . \quad (2.12)$$

We now wish to make the following points about $\text{Deff}(\hat{\beta}_{o} | \underline{x})$:

(i) When the sample sizes are all equal,

$$Deff(\hat{\beta}_{o}|\underline{x}) = 1 + (\overline{n}-1)\rho_{\chi}\rho_{\chi}$$
. (2.13)

(ii) To obtain the design effect for estimating μ under model (1.1), we let \underline{x} in (2.10) be $\frac{1}{\sqrt{n}} \frac{1}{n}$ and define the intraclass correlation for a column of 1's as 1. Then (2.13) reduces to $Deff(\hat{\mu}_0) = 1 + \left(\frac{\Sigma n_i^2}{\Sigma n_i} - 1\right) \rho_Y$, (2.14)

which also shows why
$$(2.9)$$
 is $Deff(\overline{x})$.

(iii)With equal sample sizes, (2.14) becomes

$$Deff(\hat{\mu}_{o}) = 1 + (\bar{n}-1)\rho_{\gamma}$$
, (2.15)

the well-known design effect for cluster samples.

(iv) If ρ_X and ρ_Y have the same sign, then $\text{Deff}(\beta_0|\underline{x}) > 1$, while the converse holds if ρ_X and ρ_Y have opposite signs.

(v) If either ρ_X or ρ_Y is 0, then Deff($\hat{\beta}_{\lambda}$) = 1.

(vi) If
$$\rho_{\chi} > 0$$
 and $\rho_{\chi} > 0$, then
Deff $(\hat{\beta}_{o} | \underline{x}) \le \text{Deff}(\hat{\mu}_{o})$

and

$$\operatorname{Deff}(\hat{\beta}_{x} \mid \underline{x}) \leq \operatorname{Deff}(\overline{x})$$
.

The last point is an important piece of theoretical evidence in support of Kish and Frankel's (1974) observation that design effects for complex statistics (including regression coefficients) tend to be less than design effects for means.

The fact that the design effects for means obtained from this model reduce to those used in practice for balanced samples is encouraging as is the fact that the empirical observation of Kish and Frankel is supported by the use of model (1.1).

Unfortunately, at the time of this writing, we do not have empirical values of the design effects for single variable regressions with which to compare (2.13) or (2.11). Therefore, it is not yet possible to verify the applicability of these results to sample survey situations.

3. Design Effects in a Two-Variable Regression

The model we use here is

$$\underline{y} = \underline{1}\mu + \underline{x}_1\beta_1 + \underline{x}_2\beta_2 + \underline{u}$$
 (3.1)

with
$$\underline{x}_{1}'\underline{1} = \underline{x}_{2}'\underline{1} = 0$$
, and $\underline{x}_{1}'\underline{x}_{1} = \underline{x}_{2}'\underline{x}_{2} = 1$. The variance of the OLS estimator of $(\beta_{1}, \beta_{2})'$ is

$$\operatorname{Var}\left(\begin{vmatrix} \beta_{10} \\ \beta_{20} \end{vmatrix} x\right) = (X'X)^{-1}X'VX(X'X)^{-1}, \quad (3.2)$$

where $X = [\underline{x}_1, \underline{x}_2]$. With the restriction

the restriction
$$\underline{\mathbf{x}}_1 \underline{\mathbf{x}}_1 = \underline{\mathbf{x}}_2 \underline{\mathbf{x}}_2 = 1$$
,
 $\mathbf{x}'\mathbf{x} = \begin{bmatrix} 1 & \mathbf{r} \\ \mathbf{r} & 1 \end{bmatrix}$
(3.3)

and

$$(X'X)^{-1} = \frac{1}{1-r^2} \begin{bmatrix} 1 & -r \\ -r & 1 \end{bmatrix}$$
 (3.4)

The center matrix, X'VX, is

$$\sum_{x' \vee x = (1-\rho_y) \times x' + \rho_y} \begin{bmatrix} \sum_{n_1}^2 \overline{x_{1i}}^2 & \sum_{n_1}^2 \overline{x_{1i}}^2 \overline{x_{2i}} \\ \sum_{n_1}^2 \overline{x_{1i}}^2 \overline{x_{2i}} & \sum_{n_1}^2 \overline{x_{2i}}^2 \end{bmatrix} . (3.5)$$

Using (2.9), the diagonal elements of (3.5) can be easily represented in terms of the intraclass correlations ρ_{X_1} and ρ_{X_2} . By generalizing (2.8), we define the intraclass "co-correlation" of $\underline{x_1}$ and $\underline{x_2}$ as

which reduces to

$$\rho_{\mathbf{X}_{1}\mathbf{X}_{2}} = \frac{\sum \mathbf{n_{i}^{2} \overline{\mathbf{x}_{1i} \mathbf{x}_{2i}} - \mathbf{r}}}{\sum \mathbf{n_{i}^{2}}} . \qquad (3.7)$$

We note that the sign of $\rho_{X_1X_2}$ does not depend on the sign of the covariance between the cluster totals, but on whether this covariance is larger or smaller than the overall correlation between $\underline{x_1}$ and $\underline{x_2}$. If the cluster totals are uncorrelated, then

$$\rho_{\mathbf{X}_{1}\mathbf{X}_{2}} = \frac{-\mathbf{r}}{\frac{\sum n_{i}^{2}}{\sum n_{i}} - 1}$$

To form the design effect for $\hat{\beta}_{10}$, we perform the necessary matrix multiplication in (3.2), using (2.9) and (3.7) in (3.5), to find

$$Deff(\hat{\beta}_{0}|X) = 1 + \left[\frac{\sum n_{i}^{2} \overline{x}_{1i}^{2} - 2r \sum n_{i}^{2} \overline{x}_{1i} \overline{x}_{2i}^{2} + r^{2} \sum n_{i}^{2} \overline{x}_{2i}^{2}}{1 - r^{2}} - 1\right] \rho_{Y}$$

$$= 1 + \left(\frac{\sum n_{i}^{2}}{\sum n_{i}} - 1\right) \left(\frac{\rho_{X_{1}}^{2} - 2r \rho_{X_{1}} X_{2}^{2} + r^{2} \rho_{X_{2}}}{1 - r^{2}}\right) \rho_{Y}.$$
(3.8)
(3.9)

Due to the number of parameters involved, it is difficult to make general statements about the value of $\text{Deff}(\hat{\beta}_{10}|X)$. However, we can notice that:

- (i) if r = 0, then (3.9) reduces to the single variable design effect of (2.10);
- (ii) if $\rho_{\gamma} > 0$, then Deff($\beta_{10} | X$) increases with

$$x_1$$
 and ρ_{X_2} ;

- (iii)Deff($\hat{\beta}_{10} | X$) is larger if r and $\rho_{X_1X_2}$ have opposite signs than if they have the same sign.
- (iv) $\text{Deff}(\hat{\beta}_{10}|X)$ becomes very large if r approaches 1 or -1° .

Perhaps a more intuitive parametrization of Deff($\hat{\beta}_{10} | X$) occurs when it is expressed in terms of the design effects of x_1 and x_2 . By letting ρ_{12} be the correlation coefficient between the block totals, we obtain

$$Deff(\hat{\beta}_{10}|\mathbf{X}) = (3.10)$$

$$1 + \left[\frac{Deff(\overline{\mathbf{x}}_{1}) - 2r\rho_{12}\sqrt{Deff(\overline{\mathbf{x}}_{1})Deff(\overline{\mathbf{x}}_{2})} + r^{2}Deff(\overline{\mathbf{x}}_{2})}{1 - r^{2}} - 1 \right] \rho_{\mathbf{Y}}.$$

To evaluate (3.10), we used data from Frankel's (1971) three variable regressions, considering the variables pairwise. Values of $\sqrt{\text{Deff}(\overline{x}_i)}$ were included in his appendix E. Values of r were available from Table 5.1. Values for ρ_{Y} were obtained by assuming $\text{Deff}(\overline{y}) = 1 + (\overline{n}-1)\rho_{Y}$ and solving for ρ_{Y} . As sufficient data were not available for evaluating ρ_{12} , we assumed it was equal to r. Only data from the six strata designs are used.

Frankel considered 2 different three variable regressions. Since we used the variables in pairs in 2 variable equations, (3.10) was evaluated twice for each regression coefficient. The results of our calculations and Frankel's empirically obtained values are given below.

	Comparison of	E Theoretical a	nd
	Empirical	Design Effects	
Vari- able	$\frac{\text{Deff}(\hat{\beta}_{io})}{1}$	from (3.10)	Deff(β _{io}) <u>from Frankel</u>
6	1.067,	1.063	1.089
7	1.092,	1.088	1.134
12	1.089,	1.089	.984
8	1.058,	1.057	1.093
11	1.126,	1.128	1.080
17	1,251,	1,252	1.432

Variables 6, 7, and 12 were predictor variables in one equation while 8, 11, and 17 were included in the other equation. Except for variables 12 and 11, the results from (3.10) are somewhat smaller but ordered approximately the same as Frankel's results. Variable 12 is clearly an anomaly for which we have no explanation at this time. From Frankel's data we found variables 6 and 7 were highly correlated with each other and correlated only slightly with variable 12. Until results are obtained for three variable regressions, we do not know whether this explains the small design effect for variable 12.

In this section and the preceding one, we have presented expressions for conditional design effects for regression coefficients. These were obtained by assuming the data follow a simple linear model appropriate for two-stage sampling from infinite populations. It is hoped that these results may also shed some light on the properties of regression coefficients obtained from finite populations.

The comparisons in the above table are not totally discouraging. Further investigation is needed to determine whether the discrepancies are due to differences between two-variable and three-variable regressions or from some oversimplification in the assumed model.

4. <u>Relative Efficiency of OLS for Cluster</u> <u>Samples</u>

In this section we study the efficiency of OLS with respect to WLS when model (1.1) holds. We consider only single variable regressions and define the relative efficiency as

$$E^{\star} = \frac{\operatorname{Var}(\hat{\beta}_{w}|\underline{x})}{\operatorname{Var}(\hat{\beta}_{o}|\underline{x})}, \qquad (4.1)$$

where $\hat{\boldsymbol{\beta}}_{_{\boldsymbol{W}}}$ is the second element of

$$\begin{pmatrix} \hat{\mu}_{w} \\ \hat{\beta}_{w} \end{pmatrix} = (z'v^{-1}z)^{-1}z'v^{-1}y$$

with $Z = [\underline{1}, \underline{x}]$. $Var(\hat{\beta}_w)$ is the (2,2) element of $(Z'V^{-1}Z)^{-1}$. As before we assume $\underline{x'1} = 0$ and $\underline{x'x} = 1$. The efficiencies given here are a pessimistic reflection of the efficiency of OLS since $Var(\hat{\beta}_w)$ can never be achieved.

It can be shown that

$$E^{*} = \frac{(1-\rho_{Y})}{[1-\underline{a}_{2}'D\underline{a}_{2}][1+(\underline{a}_{2}'N\underline{a}_{2}-1)\rho_{Y}]}, \quad (4.2)$$

where $\underline{a}'_{2} = [\sqrt{n_{1}}\overline{x}_{1}, \dots, \sqrt{n_{b}}\overline{x}_{b}]$

$$N = \operatorname{diag} \left(\frac{n_1, \dots, n_b}{1 + (n_1 - 1)\rho_Y}, \dots, \frac{n_b \rho_Y}{1 + (n_b - 1)\rho_Y} \right).$$

We note that

$$0 \leq \underline{a'_{2a_2}} = \Sigma n_{\underline{i}} \overline{x_{\underline{i}}}^2 \leq 1$$

is the between-cluster sum of squares of \underline{x} , since $\overline{x} = 0$. As such it is the length of the projection of \underline{x} into the subspace spanned by \underline{x}_b . The vector \underline{a}_2 contains the between-cluster information for regressing \underline{y} on \underline{x} . If $\underline{a}_2 = \underline{0}$, all cluster means are 0 and \underline{x} varies only within the clusters.

Rather than discussing the properties of E[^], we give some graphs of it in simple situations. We choose to represent \underline{a}_2 as

$$\underline{a}_2 = \sqrt{k} \underline{u}$$

where $\underline{u}'\underline{u} = 1$ and $k = \underline{a}'_2\underline{a}_2$. Using this representation, $\underline{a}'_2N\underline{a}_2 = k\underline{u}'N\underline{u}$, where the quantity $\underline{u}'N\underline{u}$ is a weighted average of the n_1 and must satisfy

$$n_1 \leq \underline{u}' N \underline{u} \leq n_b \quad . \tag{4.3}$$

The restriction $\underline{x'1} = 0$ translates to $\underline{u'\ell} = 0$

where $l' = \left(\sqrt{\frac{n_1}{n}}, \dots, \sqrt{\frac{n_b}{n}}\right)$

With this restriction on <u>u</u>, equality on the left in (4.3) can be attained only if $n_1 = n_2$ and on the right only if $n_{b-1} = n_b$.

We also point out that k is a linear function of the intraclass correlation of \underline{x} via

$$k = \frac{\rho_{X} \left(\frac{\Sigma n_{i}^{2}}{\Sigma n_{i}} - 1 \right) + 1}{\underline{u}' N \underline{u}} . \qquad (4.4)$$

For the illustrations, we consider only balanced samples and use

$$E^{*} = \frac{(1-\rho_{Y}) (1-\rho_{Y}+\overline{n}\rho_{Y})}{(1-\rho_{Y}+\overline{n}\rho_{Y}(1-k))(1-\rho_{Y}+\overline{n}\rho_{Y}k)}$$
(4.5)

which is obtained from (4.2) with $n_1 = ... = n_b = \overline{n}$. The relationship between k and ρ_X simplifies to

$$k = \frac{1 + (\bar{n} - 1)\rho_{X}}{\bar{n}} \qquad (4.6)$$

for balanced samples.

In figures 1 and 2 we present graphs of \mathbf{E}^{\star} versus k (or $\rho_{\mathbf{X}}$) for different values of $\rho_{\mathbf{Y}}$. Figure 1 contains results for $\overline{\mathbf{n}} = 100$ and Figure 2 for $\overline{\mathbf{n}} = 50$. Small values of $\rho_{\mathbf{Y}}$ and k such as are commonly found in sample survey data were used in the calculations.

If we define "reasonable efficiency" as E ≥ 0.75 , then with \overline{n} = 100 the efficiency of OLS could be unreasonably low if $\rho_Y > 0.05$ unless ρ_X is very small. With \overline{n} = 50, the values of E remain high until $\rho_Y > 0.10$. Given the small values of ρ_Y and ρ_X commonly present in social science data, OLS should be reasonably efficient for most ρ_X and ρ_Y when $\overline{n} \leq 50$. With large values of \overline{n} , some inefficient estimates may result if OLS is used consistently.

We also point out that with large total sample sizes, \overline{n} and/or b both large, then Var($\hat{\beta}$) and Var($\hat{\beta}_w$)may both be acceptably small even though the efficiency of OLS is low - making it not worthwhile to attempt a WLS analysis.

In conclusion, with the presence of clustering as modelled in (1.1), it appears that OLS is a reasonably efficient estimator of a single regression coefficient for many parameter values commonly obtained in social science data. Therefore standard computer programs can usually be used for calculating point estimates of regression coefficients. The properties of the estimated standard error of $\hat{\beta}_0$ provided by an OLS routine when model (1.1) holds have not been investigated at this time.

REFERENCES

- Frankel, Martin R. (1971) Inference from Sample Surveys: An Empirical Investigation. Ann Arbor: Institute for Social Research, The University of Michigan.
- Kish, Leslie and Frankel, Martin R. (1974). "Inference from complex samples." <u>The</u> <u>Journal of the Royal Statistical Society</u>, <u>Series B</u>, <u>36</u>. No. 1, pp. 1-37.
- Murthy, M.N. (1967). <u>Sampling: Theory and</u> <u>Methods</u>. Calcutta: Statistical Publishing Society.

ACKNOWLEDGMENT

This research is based, in part, on the author's Ph.D. dissertation which was completed at Southern Methodist University.



Figure 1 E* vs. K(ρ_x) for Different Values of ρ_y : $\bar{n} = 100$



Figure 2 E* vs. ρ_x for Different Values of ρ_y : $\bar{n} = 50$