PAUL S. LEVY University of Illinois at the Medical Center School of Public Health

1. Introduction

The use of sample surveys with multiplicity has been advocated by Sirken (1970) for estimating the number of demographic events (e.g. births, deaths) occurring in a particular time period. Since Sirken's original article, the theory of multiplicity estimation (also called network sampling) has been extended to stratified random sampling (Sirken, 1972; Levy, 1973), estimation of proportions (Sirken and Levy, 1974) and estimation of response errors (Nathan, 1976). In addition, sample surveys with multiplicity have been used in a wide variety of applications (Sirken, 1972; Sirken and Levy, 1974; Sirken et al., 1975; Nathan et al., 1977). In this report, the theory of multiplicity estimation is extended to simple cluster sampling, and an unbiased estimator is proposed for estimating the total number of events under this type of sampling design.

A survey with multiplicity is one in which an element (e.g. birth, death, individual having some attribute, etc.) may be linked to more than one enumeration unit by an algorithm or counting rule. For example, a counting rule in a survey with multiplicity might link a birth to the households of the grandparents as well as to the parents' household whereas a conventional counting rule would link the birth only to the household of the parents.

2. Development of the Estimator:

2.1 Population Parameters

Let us suppose that a population contains N enumeration units (e.u.'s) grouped into M primary sampling units (PSU's) with PSU i containing N_i e.u.'s; i = 1, ..., M, and that a counting rule links Y events labeled I₁, ..., I_Y to enumeration units according to an indicator variable, $\delta_{\alpha ij}$ given by

$$\delta_{\alpha i j} = \begin{cases} \text{if event } I_{\alpha} \text{ is linked to e.u. j} \\ \text{in PSU i by the counting rule} \\ \text{otherwise} \end{cases}$$

where

 $\alpha = 1$, ..., Y, i = 1, ..., M and j = 1, ..., N_i.

For any counting rule, the following parameters can be defined which characterize the network linking the enumeration units to the elements:

$$s_{\alpha i} = \sum_{j=1}^{N} \delta_{\alpha i j};$$

$$s_{\alpha} = \sum_{i=1}^{M} s_{\alpha i} ;$$

$$t_{\alpha i} = \begin{cases} 1 & \text{if } s_{\alpha i} > 0 \\ 0 & \text{if } s_{\alpha i} = 0 \end{cases}$$

$$t_{\alpha} = \sum_{i=1}^{M} t_{\alpha i}$$

The parameter, $S_{\alpha i}$, denotes the number of enumeration units in a particular PSU (i) that are linked to a particular element, I_{α} , whereas s_{α} denotes the total number of enumeration units linked to an element by a counting rule and is referred to as the <u>multiplicity</u> of the element with respect to the counting rule. Clearly, for conventional counting rule, each s_{α} would be equal to unity. The parameter, t_{α} , denotes the number of PSU's in which a particular element, I_{α} , is linked to one or more enumeration units.

Let $(z_{\alpha i} : \alpha = 1, ..., Y; i = 1, ..., M;$ $t_{\alpha i} = 1$) be any set of weights defined for all (α, i) such that $t_{\alpha i} = 1$ with the property:

$$\sum_{i=1}^{M} z_{\alpha i} s_{\alpha i} = 1 \qquad \alpha = 1, \dots, Y$$

We then define the following parameters for each PSU.

$$\lambda^{*}_{ij} = \sum_{\alpha=1}^{Y} Z_{\alpha i} \delta_{\alpha i j} , \quad j = 1, \dots, N_{i}$$
$$Y_{i}^{*} = \sum_{j=1}^{N_{i}} \lambda^{*}_{i j} ; \quad E_{i} = \sum_{\alpha=1}^{Y} Z_{\alpha i}^{2} S_{\alpha i} / Y_{i}^{*}$$

The parameter, λ_{ij}^{\prime} , represents the basic summary information obtained from enumeration units concerning elements, while the weights $\{z_{\alpha i}\}$ are functions of the particular network linking enumeration to elements and are chosen to make estimates of Y unbiased. Some possible choices of $z_{\alpha i}$ might be $1/s_{\alpha}$ or $1/(s_{\alpha i} t_{\alpha})$ for those counting rules which link elements to enumeration units in more than one PSU. For counting rules which link elements to enumeration units in only one PSU, the $z_{\alpha i}$ might be set equal to $1/s_{\alpha i}$, and for conventional counting rules, the $z_{\alpha i}$ would be equal to 1. The E_i are generalizations of parameters found by Sirken and Levy (1974) to be involved in the variances of estimates obtained from multiplicity surveys, while the parameters, Y_1^* , although not necessarily integer valued, could be interpreted as being the "effective number" of elements linked to PSU i by the enumeration rule. It can be shown that M \star

Using nomenclature similar to that in Hansen, Hurwitz, and Madow (1953), we define for the variable, λ'_{ij} , the between PSU variance

$$s_{1Y}^{2} * = \sum_{i=1}^{M} (Y_{i} * - \overline{Y})^{2} / (M - 1)$$

the within PSU variance,

$$s_{2Y}^{2} \star = \frac{1}{N} \frac{M}{i = 1} \frac{M}{M_{i} - 1} \frac{N_{i}}{j = 1} (\lambda'_{ij} - Y_{i}^{\star})^{2}$$

and the intra-class correlation coefficient

$$\delta = \left(\frac{M-1}{M} \quad s_{1Y}^2 \star - \bar{N} \quad s_{2Y}^2 \star \right) / \frac{(M-1)}{M} \quad s_{1Y}^2 \star + \bar{N} \quad (\bar{N}-1) \quad s_{2Y}^2 \star$$

where

$$\overline{\mathbf{Y}} = \frac{\mathbf{M}}{\sum_{i=1}^{L} \mathbf{Y}_{i}} \times \mathbf{M}$$
$$= \mathbf{Y}/\mathbf{M} \text{ (since } \mathbf{M}_{i=1}^{M} \mathbf{Y}_{i} = \mathbf{Y})$$

and

$$\bar{N} = \sum_{i=1}^{N} N_i / M$$

М

With these definitions, the following theorems can be proved.

Theorem 1.

The variance, $S_{1Y}^2 *$ among PSU's with respect to λ'_{ij} is equal to the expression given by:

$$s_{1Y}^{2} * = [M \ \overline{Y} (E_{k} - \overline{Y}) + \sum_{\substack{\alpha = 1 \ \alpha' = 1}}^{Y} \sum_{\substack{\alpha' = 1 \ \alpha' \neq \alpha}}^{Y} v_{\alpha \alpha'}] / (M-1).$$

where

$$E_{k} = \sum_{\alpha i} \sum_{\alpha i} z_{\alpha i}^{2} s_{\alpha i}^{2} / Y$$

and

$$v_{\alpha\alpha'} = \sum_{i}^{\Sigma} z_{\alpha i} \quad z_{\alpha'i} \quad s_{\alpha i} \quad s_{\alpha'i} \quad ; \quad \alpha' \neq \alpha.$$

PROOF

$$\sum_{i=1}^{M} (Y_{i}^{*} - \bar{Y})^{2} = \sum_{i=1}^{M} (Y_{i}^{*})^{2} - M \bar{Y}^{2}$$

$$= \sum_{i=1}^{M} (\sum_{j=1}^{N} \sum_{\alpha=1}^{Y} z_{\alpha i} \delta_{\alpha i j})^{2} - M \overline{Y}^{2}$$

$$= \sum_{i=1}^{M} (\sum_{\alpha=1}^{Y} z_{\alpha i} s_{\alpha i})^{2} - M \overline{Y}^{2}$$

$$= \sum_{i=1}^{M} \sum_{\alpha=1}^{Y} z_{\alpha i}^{2} s_{\alpha i}^{2} + \sum_{i=\alpha=1}^{Y} \sum_{\alpha=1}^{Y} z_{\alpha i} z_{\alpha'} s_{\alpha i}^{\beta} s_{\alpha' i}$$

$$- M \overline{Y}^{2}$$

$$= M \overline{Y} (E_{k} - \overline{Y}) + \sum_{\alpha=1}^{Y} \alpha' z_{\alpha'}$$

$$\alpha' \neq \alpha$$

$$q.e.d.$$

Theorem 2.

The within PSU variance with respect to $\lambda_{\mbox{ij}}^{\mbox{!}}$ is given by

$$s_{2Y*}^{2} = \frac{1}{N} \prod_{i=1}^{M} \frac{N_{i}^{2}}{N_{i}^{-1}} \overline{Y}_{i}^{*} (E_{i} - \overline{Y}_{i}^{*})$$
$$+ \frac{1}{N} \prod_{i=1}^{M} \frac{N_{i}}{N_{i}^{-1}} \alpha_{\alpha}^{Y} \sum_{\alpha \in 1}^{Y} \gamma_{i\alpha\alpha}^{*}$$
$$\alpha^{*} \neq \alpha$$

where
$$v_{i\alpha\alpha'} = \frac{N_i}{\sum_{i=1}^{\Sigma} z_{\alpha i} z_{\alpha' i} \delta_{\alpha i j} \delta_{\alpha' i j}} for \alpha' \neq \alpha$$

$$\begin{split} & \underset{i \stackrel{M}{\leq} 1}{\overset{N_{i}}{\underset{N_{i}^{-1}}{\overset{N_{i}}{=} 1}} } \sum_{j \stackrel{M}{\leq} 1}^{\overset{N_{i}}{=} 1} (\lambda_{ij}^{*} - \overline{y}_{i}^{*})^{2} \\ &= \underset{i \stackrel{M}{\leq} 1}{\overset{N_{i}^{-1}}{\underset{N_{i}^{-1}}{\overset{N_{i}^{-1}}{\underset{j \stackrel{\Sigma}{=} 1}{\overset{N_{i}}{\underset{(j \stackrel{\Sigma}{=} 1}{\overset{N_{i}}{\underset{(j \stackrel{\Sigma}{=} 1}{\overset{X_{ij}}{\underset{(j \stackrel{\Sigma}{=} 1}{\overset{X_{ij}}{\underset{(j \stackrel{\Sigma}{=} 1}{\overset{X_{ij}}{\underset{(j \stackrel{\Sigma}{=} 1}{\underset{(j \stackrel{\Sigma}{=} 1}{\underset{$$

$$= \underbrace{\prod_{i=1}^{M} \frac{N_{i}}{N_{i}-1} \left(\sum_{\alpha=1}^{Y} z_{\alpha i}^{2} s_{\alpha i} + \sum_{i=1}^{N} \sum_{\alpha=1}^{Y} \frac{Y}{\alpha} \right)}_{z_{\alpha i} z_{\alpha' i} \delta_{\alpha i j} \delta_{\alpha' i j} - N_{i} \left(\overline{Y}_{i}^{*}\right)^{2}}$$

$$= \sum_{i=1}^{M} \frac{N_{i}}{N_{i}-1} (N_{i} \overline{Y}_{i}^{*} (E_{i} - \overline{Y}_{i}^{*}) + \sum_{\alpha \equiv 1 \alpha}^{Y} \overline{\Sigma}_{1} \nabla_{i\alpha\alpha})$$

$$\alpha' \neq \alpha$$

Corollary 1.

The intra class correlation coefficient, $\delta,$ is given by

$$\delta = \frac{A - B}{A + (N-1)B}$$

where

$$A = M \overline{Y} (E_{k} - \overline{Y}) + \sum_{\alpha \in I} \sum_{\alpha' \neq \alpha} v_{\alpha' \neq \alpha}$$

and

$$B = \sum_{i=1}^{N} \frac{N_L^2}{N_i - 1} \overline{Y}_i^* (E_i - \overline{Y}_i^*)$$
$$+ \sum_{i=1}^{M} \frac{N_L^2 Y}{N_i - 1\alpha} \sum_{\substack{i=1 \\ \alpha \neq \alpha}} V_{i\alpha\alpha}$$

Proof follows from Theorems 1 and 2 and the definition of δ .

q.e.d.

Corollary 2.

ŝ

If the assumption is made that an enumeration unit is linked to no more than one element then $v_{i\alpha\alpha'} = 0$ for all i, α , and α' ; s_{2Y}^{2} is given by

$$S_{2Y}^{2} \star = \frac{1}{N} \sum_{i=1}^{M} \frac{N_{i}^{2}}{N_{i}-1} \quad (E_{i} - \overline{Y}_{i}^{\star})$$

and

$$\delta = \frac{A - B}{A + (N-1)B}$$

where

$$A = M \overline{Y} (\underline{E}_{k} - \overline{Y}) + \sum_{\alpha=1}^{Y} \alpha, \sum_{\alpha} v_{\alpha\alpha},$$

α" ≠ α

and

$$B = \frac{M}{1} \frac{N_{1}^{2}}{N_{1}^{2} - 1} \frac{\overline{Y}_{1}^{*}}{\overline{Y}_{1}^{*}} (E_{1}^{*} - \overline{Y}_{1}^{*})$$

2.2 Estimation of Y, the Total Number of Events from the Sample

Let us assume that the sample design used to estimate the number of events, Y, in the population is a simple two stage cluster sample as defined by Hansen, Hurwitz, and Madow (1953). In other words, a simple random sample of m PSU's is taken from the M PSU's in the population, and within each sample PSU (i), a simple random sample of n_i enumeration units is taken from the N_i enumeration units in the PSU, with the second stage sampling fraction, n_i/N_i , the same for each PSU.

If (for convenience) the sample PSU's are labelled 1,, m and the sample enumeration units within each sample PSU are denoted i₁, ..., i_n where i = 1, ..., m, then the estii mator Y' of Y as given by

$$\mathbf{Y}' = \frac{\mathbf{M}}{\mathbf{m}} \sum_{i=1}^{m} \frac{\mathbf{N}_{i}}{\mathbf{n}_{i}} \sum_{j=1}^{n} \lambda_{ij}' \quad (1)$$

is an unbiased estimator of Y as shown below in Theorem 3.

Theorem 3.

The estimator, Y' of Y as defined in equation (1) is an unbiased estimator of Y.

PROOF

For a given sample PSU, i, the expected value over all possible second stage samples of n.

$$\sum_{j=1}^{j} \lambda^{i}$$
 is given by

$$E(\sum_{j=1}^{n_{i}} \lambda_{ij}^{\prime}) = \frac{n_{i}}{n_{i}} \sum_{j=1}^{n_{i}} \lambda_{ij}^{\prime} = \frac{n_{i}}{n_{i}} Y_{i}^{\prime}$$

Thus, the expected value of Y' over all possible samples is given by

$$E(Y^*) = \stackrel{M}{m} \stackrel{E}{all i} \stackrel{M}{\underset{i=1}{\overset{M}{i=1}}} \stackrel{M}{\underset{n_i}{\overset{N_i}{i=1}}} \stackrel{E}{\underset{n_i}{\overset{n_i}{i=1}}} \stackrel{n_i}{\underset{n_i}{\overset{N_i}{i=1}}} \stackrel{n_i}{\underset{i=1}{\overset{N_i}{j=1}}} \stackrel{n_i}{\underset{i=1}{\overset{N_i}{\underset{i=1}{\overset{N_i}{\underset{i=1}{\underset{i=1}{\overset{N_i}{j=1}}}} \stackrel{n_i}{\underset{i=1}{\underset{i=1}{\overset{N_i}{\underset{i=1}{\underset{i=1}{\underset{i=1}{\overset{N_i}{j=1}}}} \stackrel{n_i}{\underset{i=1$$

3. Some Relationships involving
$$\delta$$
 When all $v_{igg} = 0$

When $V_{i\alpha\alpha}$ = 0 for all $\alpha_i \alpha$ ' and for all i, then the intraclass correlation coefficient, δ_i , is given by $\delta = (A - B) / (A + (\overline{N} - 1) B)$

where

 $A = M \overline{Y} (E_{k} - \overline{Y}) + \Sigma \Sigma V$ $\alpha \alpha' \alpha \alpha'$ $\alpha' \neq \alpha$

and

 $B = \sum_{i=1}^{\Sigma} N_i \tilde{Y}_i^* (E_i - \tilde{Y}_i^*) .$ Clearly $A \ge 0$ and $B \ge 0$ since they are quadra-

$$\frac{d \delta}{d A} = \frac{B}{\left[A + (\overline{N} - 1) B\right]^2} \ge 0 ,$$

and therefore, δ varies directly with A.

tic forms. It can be shown that since B > 0

On the other hand, since A > 0, then

$$\frac{d \delta}{d B} = \frac{-\overline{N} A}{[A + (\overline{N} - 1) B]} 2 \stackrel{>}{=} 0$$

and hence δ varies inversely with B.

Let us examine δ for the set of weights $z_{\alpha i} = 1/(t_{\alpha}s_{\alpha i})$; $\alpha = 1, \ldots, Y$; $i = 1, \ldots, M$. For this set of weights, we have:



$$\bar{\mathbf{Y}}_{i}^{\star} = \frac{1}{N-1} \begin{array}{c} \mathbf{Y} & \mathbf{t} \\ \boldsymbol{\Sigma} & \boldsymbol{\alpha}^{i} \\ \boldsymbol{\alpha}^{=1} & \mathbf{t} \\ \boldsymbol{\alpha} \end{array}$$

 $E_{k} = \frac{Y}{\alpha=1} \frac{1}{t_{\alpha}} / Y$

and

If; the multiplicities, $s_{\alpha i}$, are increased without increasing the t or t , then the α E, would decrease which would cause a decrease in B since the \bar{Y}_{i}^{\star} would be unaffected. Since, also, the $v_{\alpha\alpha'}$ and E_{k} would not be affected by change in the $s_{\alpha i}$, it follows that A would not be affected. Hence, increase in sai would result in an increase in the intra-class correlation coefficient, δ .

REFERENCES

Sirken, M.G. (1970): "Household Surveys

with Multiplicity" Journal of the American Statistical Association 65, 257 - 266.

2. Sirken, M.G. (1972): "Stratified Sample Surveys with Multiplicity" <u>Journal of the</u> American Statistical Association 67, 224 - 227.

Levy, P. S. (1978): "Optimum Allocation in Stratified Random Network Sampling for Estimating the Prevalence of Attributes in Rare Populations" To appear in Journal of the American Statistical Association, December, 1977.

Sirken, M.G. and Levy, P.S. (1974): "Multiplicity Estimation of Proportions Based on Ratios of Random Variables" Journal of the American Statistical Association 69, 68 - 73.

5. Sirken, M. G., Indurfurth, G. P., Burnham, C. E., and Danchik, K. M. (1975): "Household Sample Surveys of Diabetes: Design Effects of Counting Rules" American Statistical Association, Proceedings of the Social Statistics Section, 659 - 663.

Nathan, G., Schmelz, U. O. and Kenvin, J. 6. (1977): "Multiplicity Study of Marriages and Births in Israel" Vital and Health Statistics, Series 2 No. 70, NCHS, Rockville, Maryland.

Nathan, G. (1976): "An Empirical Study 7. of Response and Sampling Errors for Multiplicity Estimates with Different Counting Rules" Journal of the American Statistical Association 71, 808 - 815.