

Aridaman K. Jain, Bell Telephone Laboratories

1. INTRODUCTION

The problem of changing strata and the resultant changes in selection probabilities for primary units in panel studies has been discussed by Kish [1], Fellegi [2,3] and others [4,5,6,7,8]. After the initial selection, the sampled units may be used for several surveys over a period of several years. During this period, there may be significant and frequent changes in the population. Since changes would continue to occur in the future, a single brand new selection of sample units is not a satisfactory solution. Moreover, continued use of the original sample has several advantages over a brand new sample: (i) it is much cheaper, (ii) it provides a more precise comparison of periodic results since the initial selection, and (iii) it avoids the delay in the availability of results due to the start-up of a new sample.

We will consider the case when units within a stratum are selected with probability proportional to initial size (pps) and the value of one of the stratification factors may change over time. Information on such changes, which result in units moving from their original strata to new strata, is readily available for the sampled units but not for the nonsampled units. As a result of changes in the composition of strata, the stratified sample no longer corresponds to the original pps selection within each stratum. A method of estimation based on the original selection of the sample but incorporating the subsequent changes in the sample is proposed. Major advantages of the proposed procedure, which yields approximately unbiased estimates, are that it does not require

- (i) information on changes in the nonsampled units,
- (ii) replacement of some units in the original sample by new units.

Here we consider panels of plant facilities where type of facility is a principal factor for stratification. But, the techniques discussed here apply to panel studies in general. For example, a consumer panel may be a stratified random sample with family income as a stratification factor. Then, a significant change in family income of a panel member would correspond to a change of type of facility.

2. DESIGN OF THE PANEL AND NOTATION

Two stratification factors, namely, geographical area and type of switching machine, were used to stratify the population of originating entities in a state. Within each stratum a sample of entities was chosen with replacement with probabilities proportional to initial size (number of working lines). Then several experimental lines were assigned to each sample entity. Several types of calls are made from these experimental lines on a continuing basis. The details of these calls are recorded by a mini-

computer to determine the number of recording errors.

Let

i = geographical area, $i = 1, 2, \dots, 7$;

j = original switching machine type, $j = 1, 2, \dots, 5$;

u = current switching machine type, $u = 1, 2, \dots, 5$;

t = call type, t = single message unit, multiple message unit, toll, operator handled;

M_{ij} = total number of originating entities in (i, j) ;

m_{ij} = sample number of distinct originating entities in (i, j) ;

w_{ijq} = number of times entity q is in the sample, $q = 1, 2, \dots, m_{ij}$;

m_{iju} = number of distinct sample entities in (i, j) out of m_{ij} which are currently of switch type u , where $\sum_{u=1}^5 m_{iju} = m_{ij}$;

v_{ijq} = initial number of working lines in entity q in (i, j) ;

v_{ijq}^N = current number of working lines in entity (i, j, q) ;

c_{tijq} = call volume of type t for (i, j, q) ;
 c_{tijqu} = c_{tijq} if the current switch type of entity (i, j, q) is u , 0 otherwise;

l_{ijq} = number of experimental lines assigned to (i, j, q) ;

\bar{n}_{tijq} = average number of experimental calls of type t per experimental line in (i, j, q) ;

P_{tijqr} = probability of error for call type t for experimental line r in (i, j, q) , $r = 1, 2, \dots, l_{ijq}$;

s = call number;

x_{tijqrs} = 1 if a test call is the s th call of type t originating from experimental line r in entity (i, j, q) , 0 otherwise;

y_{tijqrs} = 1 if $x_{tijqrs} = 1$ and that test call is in error, 0 otherwise.

As described above, the panel of facilities consists of a sample of entities for each switching machine type for each geographical area. After the initial selection, the process of modernization results in the replacement of entities of older types by entities of newer types. Such a replacement results in the movement of a panel entity from one stratum to another stratum and consequently it changes strata composition and selection probabilities. The next section describes an estimation procedure which does not require information on changes in the nonsampled entities.

3. ESTIMATION PROCEDURE

Because of the replacement of some entities, our interest has shifted from original strata (i,j), from which the panel members were selected, to new strata (i,u) which correspond to the current breakdown of the population of entities. The estimates for new strata (i,u) are derived by first subdividing each original stratum (i,j) into domains of current switch type (u=1,2,...,5) and then summing over j for each geographical area i.

3.1 Estimates for an Entity in (i,j)

An estimate of p_{tijq} , the probability of error for call type t for entity (i,j,q), is given by the following conventional ratio estimator:

$$\hat{p}_{tijq} = \frac{\hat{Y}_{tijq}}{\hat{X}_{tijq}}, \quad (1)$$

where

$$\hat{X}_{tijq} = \left(\frac{v_{ijq}^N}{l_{ijq}} \right) \sum_r \hat{X}_{tijqr},$$

$$\hat{Y}_{tijq} = \left(\frac{v_{ijq}^N}{l_{ijq}} \right) \sum_r \hat{Y}_{tijqr},$$

$$\hat{X}_{tijqr} = \sum_s x_{tijqrs},$$

$$\hat{Y}_{tijqr} = \sum_s y_{tijqrs}.$$

The estimator \hat{p}_{tijq} can be rewritten as

$$\hat{p}_{tijq} = \sum_r \left(\frac{\hat{X}_{tijqr}}{\sum_r \hat{X}_{tijqr}} \right) \hat{p}_{tijqr}$$

where

$$\hat{p}_{tijqr} = \frac{\hat{Y}_{tijqr}}{\hat{X}_{tijqr}}.$$

It can be shown that

$$E(\hat{p}_{tijqr}) = p_{tijqr}$$

and

$$E(\hat{p}_{tijq}) = \bar{p}_{tijq}.$$

where \bar{p}_{tijq} is a weighted average of line probabilities.

It may be noted that \hat{X}_{tijq} is an estimate of total customer call volume of type t for entity (i,j,q) under the assumption that on an average customers make n_{tijq} calls per line. In subsequent discussion we will replace a subscript by

a dot to indicate either summation over test calls (X and Y) or averaging for probability of error (p). An estimate of the variance of \hat{p}_{tijq} is given by

$$v(\hat{p}_{tijq}) = \frac{\sum_{r=1}^{l_{ijq}} [\hat{Y}_{tijqr} - \hat{p}_{tijq} \hat{X}_{tijqr}]^2}{l_{ijq} (l_{ijq} - 1) (\bar{n}_{tijq})^2}. \quad (2)$$

3.2 Estimate for Domain u in (i,j)

In the subsequent discussion, since the subscript in the 5th place is not needed to denote the line number (we have already summed over lines), it will be used to denote the current switch type (u). If $m_{iju} = 0$, we define

$$\hat{X}_{tij.u} = \hat{Y}_{tij.u} = \hat{p}_{tij.u} = v(\hat{p}_{tij.u}) = 0. \quad (3)$$

If $m_{iju} \geq 1$, we define

$$\hat{X}_{tij.u} = \frac{1}{w_{ij.}} \sum_{q=1}^{m_{ij}} \left(\frac{v_{ij.}}{v_{ijq}} \right) (c_{tijqu}) (w_{ijq}) \quad (4)$$

$$\hat{Y}_{tij.u} = \frac{1}{w_{ij.}} \sum_{q=1}^{m_{ij}} \left(\frac{v_{ij.}}{v_{ijq}} \right) (c_{tijqu} \cdot \hat{p}_{tijq}) (w_{ijq}) \quad (5)$$

$$\hat{p}_{tij.u} = \hat{Y}_{tij.u} / \hat{X}_{tij.u} \quad (6)$$

and

$$v(\hat{p}_{tij.u}) = \frac{1}{w_{ij.} (w_{ij.} - 1) (\hat{X}_{tij.u})^2} \sum_{q=1}^{m_{ij}} \left[\left(\frac{v_{ij.}}{v_{ijq}} \right) (c_{tijqu}) \right]^2 (w_{ijq}) \cdot (\hat{p}_{ijq} - \hat{p}_{tij.u})^2 \quad (7)$$

where

$$v_{ij.} = \sum_{q=1}^{M_{ij}} v_{ijq} \quad \text{and} \quad w_{ij.} = \sum_{q=1}^{m_{ij}} w_{ijq}.$$

If $m_{iju} = 1$, then there exists a value of q (say q_1) such that $p_{tij.u} = p_{tijq_1}$ and Formula (7) is not applicable. In this case we define

$$v(\hat{p}_{tij.u}) = v(\hat{p}_{tijq_1}). \quad (8)$$

Approximate unbiasedness of the estimator $\hat{p}_{tij.u}$ is shown in Appendix A.

3.3 Estimates for (i,u)

Now for each combination of domain (u) and geographical area (i) we combine the five strata corresponding to the original switch type (j),

$$\hat{X}_{ti\cdot\cdot u} = \sum_j \hat{X}_{tij\cdot u} \quad (9)$$

$$\hat{Y}_{ti\cdot\cdot u} = \sum_j \hat{Y}_{tij\cdot u} = \sum_j \hat{p}_{tij\cdot u} \hat{X}_{tij\cdot u} \quad (10)$$

$$\hat{p}_{ti\cdot\cdot u} = \hat{Y}_{ti\cdot\cdot u} / \hat{X}_{ti\cdot\cdot u} = \sum_j \frac{\hat{X}_{tij\cdot u}}{\hat{X}_{ti\cdot\cdot u}} \hat{p}_{tij\cdot u} \quad (11)$$

$$v(\hat{p}_{ti\cdot\cdot u}) = \sum_j \left[\frac{\hat{X}_{tij\cdot u}}{\hat{X}_{ti\cdot\cdot u}} \right]^2 v(\hat{p}_{tij\cdot u}) \quad (12)$$

It may be noted that the above variance formula is an approximation (underestimate) due to the fact that the variation in the ratios $\hat{X}_{tij\cdot u} / \hat{X}_{ti\cdot\cdot u}$ is being ignored. A better approximation is given in Appendix B. It is assumed in the above that $\hat{X}_{ti\cdot\cdot u}$ is not zero. If it were, then $\hat{p}_{ti\cdot\cdot u}$ is not estimable.

It may also be noted that if none of the sample originating entities is replaced by one of type u, then

$$\hat{p}_{ti\cdot\cdot u} = \hat{p}_{tiu\cdot u} \text{ and } v(\hat{p}_{ti\cdot\cdot u}) = v(\hat{p}_{tiu\cdot u})$$

since $\hat{X}_{ti\cdot\cdot u} = \hat{X}_{tiu\cdot u}$. In this case the above variance formula is exact.

4. SUMMARY AND GENERALIZATIONS

In panel studies, the sampled units are used for several surveys over a period of years. The problem of changing strata and the resultant changes in selection probabilities after the initial selection is not solved by a single brand new selection of sampled units. Here we have discussed a method of estimation based on the original selection of the sample but incorporating subsequent changes in the value of one of the stratification factors. This method does not require (i) information on changes in the non-sampled units, and (ii) replacement of some units in the original sample by new units.

The estimation technique described above assumes that the modernization process does not result in the rearrangement of sampling units (i.e., a sampling unit stays intact, only its type is changed). One generalization would be to the case of more complex changes (e.g., three units rearranged into two units).

Another generalization would be to develop a scheme for periodic updates so that only information on changes since the last update would be required. The current estimation technique requires the knowledge of all changes in sampled units since the original selection.

REFERENCES

- [1] Kish, L., "Changing Strata and Selection Probabilities," Proceedings of the Social Statistics Section, American Statistical Association, Washington, D.C., (1963), 124-131.
- [2] Fellegi, I. P., "Changing the Probabilities of Selection When Two Units are Selected with PPS without Replacement," Proceedings of the Social Statistics Section, American Statistical Association, Washington, D.C., (1976), 434-442.
- [3] Fellegi, I. P., "Sampling with Varying Probabilities without Replacement: Rotating and Non-rotating Samples," JASA, 58 (1963), 183-201.
- [4] Hansen, M. H. Hurwitz, W. N., and Madow, W.G., "Sample Survey Methods and Theory," Volume 1, Wiley (1953).
- [5] Keyfitz, N., "Sampling with Probabilities Proportional to Size: Adjustment for Changes in the Probabilities," JASA, 46 (1951), 105-109.
- [6] Kish, L. and Hess, I., "Some Sampling Techniques for Continuing Survey Operations," Proceedings of the Social Statistics Section, American Statistical Association, Washington, D.C., (1959), 139-143.
- [7] Kish, L. and Scott, A., "Retaining Units After Changing Strata and Probabilities," JASA, 66 (1971), 461-470.
- [8] Platek, R. and Singh, M. P., "A Strategy for Updating Continuous Surveys," Survey Methodology (Statistical Services, Statistics Canada), Volume 1 (1975), 16-26.

APPENDIX A

PROOF OF APPROXIMATE UNBIASEDNESS OF $\hat{p}_{tij\cdot u}$

There are m_{iju} distinct entities in the sample in (i,j) which are currently of switch type u. Let these be entity 1,2,..., m_{iju} in cell (i,j).

Similarly, let there be M_{ij} entities in the population in cell (i,j). Of these let M_{iju} be currently of switch type u. Then the conditional probability of selection of an entity in (i,j,u)

is $v_{ijqu} / v_{ij\cdot u}$ where $v_{ij\cdot u} = \sum_{q=1}^{M_{iju}} v_{ijqu}$ and v_{ijqu}

is the size of entity q in (i,j) which is currently of switch type u. The estimated call volume* in M_{iju} entities of type u is

$$\hat{X}_{tij \cdot u} = \left(\sum_{q=1}^{m_{iju}} \left\{ (c_{tijq}) / \left[\frac{v_{ijqu}}{v_{ij \cdot u}} \right] w_{ijq} \right\} \right) / \left(\sum_{q=1}^{m_{iju}} w_{ijq} \right)$$

But, it is quite difficult to determine $v_{ij \cdot u}$.

Note that Probability (an entity of type u is selected) = $v_{ij \cdot u} / v_{ij \cdot}$. Since m_{iju} of the m_{ij} sample entities are of type u, an unbiased

estimate of $\frac{v_{ij \cdot u}}{v_{ij \cdot}}$ is $\frac{\sum_{q=1}^{m_{iju}} w_{ijq}}{\sum_{q=1}^{m_{ij}} w_{ijq}}$.

$$\begin{aligned} \therefore \hat{X}_{tij \cdot u} &= \left(\sum_{q=1}^{m_{iju}} w_{ijq} \cdot c_{tijq} \cdot \frac{v_{ij \cdot}}{v_{ijqu}} \right) \\ &\cdot \left(\frac{\sum_{q=1}^{m_{iju}} w_{ijq}}{\sum_{q=1}^{m_{ij}} w_{ijq}} \right) / \sum_{q=1}^{m_{iju}} w_{ijq} \\ &= \frac{1}{w_{ij \cdot}} \sum_{q=1}^{m_{iju}} w_{ijq} \cdot c_{tijq} \cdot \frac{v_{ij \cdot}}{v_{ijqu}} \\ &= \frac{1}{w_{ij \cdot}} \sum_{q=1}^{m_{ij}} w_{ijq} c_{tijqu} \cdot \frac{v_{ij \cdot}}{v_{ijqu}} \end{aligned}$$

$$\therefore E[\hat{X}_{tij \cdot u}] = E \left(\frac{1}{w_{ij \cdot}} \sum_{q=1}^{m_{iju}} w_{ijq} \cdot c_{tijq} \cdot \frac{v_{ij \cdot}}{v_{ijqu}} \right)$$

We take this expected value in two steps. First, we compute the expected value conditional on m_{iju} and then over m_{iju} .

$$\begin{aligned} E(\hat{X}_{tij \cdot u} | m_{iju}) &= \frac{v_{ij \cdot}}{v_{ij \cdot u}} \cdot \left(\frac{\sum_i^{m_{iju}} w_{ijq}}{w_{ij \cdot}} \right) \\ &\cdot E \left(c_{tijq} \cdot \frac{v_{ij \cdot u}}{v_{ijqu}} \right) \\ &= \frac{v_{ij \cdot}}{v_{ij \cdot u}} \cdot \frac{\sum_i^{m_{iju}} w_{ijq}}{w_{ij \cdot}} \cdot c_{tij \cdot u} \end{aligned}$$

where $c_{tij \cdot u}$ = total call volume in M_{iju} entities of type u. Now taking the expected values over m_{iju} ,

$$\begin{aligned} E(\hat{X}_{tij \cdot u}) &= \frac{v_{ij \cdot}}{v_{ij \cdot u}} \cdot \frac{v_{ij \cdot u}}{v_{ij \cdot}} \cdot c_{tij \cdot u} \\ &= c_{tij \cdot u} = \sum_{q=1}^{M_{iju}} c_{tijq} = X_{tij \cdot u} \end{aligned}$$

Similarly,

$$\begin{aligned} \hat{Y}_{tij \cdot u} &= \frac{1}{w_{ij \cdot}} \sum_{q=1}^{m_{iju}} w_{ijq} \cdot c_{tijq} \cdot \hat{p}_{tijq} \cdot \frac{v_{ij \cdot}}{v_{ijqu}} \\ &= \frac{1}{w_{ij \cdot}} \sum_{q=1}^{m_{ij}} w_{ijq} \cdot c_{tijqu} \cdot \hat{p}_{tijq} \cdot \frac{v_{ij \cdot}}{v_{ijqu}} \end{aligned}$$

Now,

$$\begin{aligned} E(\hat{Y}_{tij \cdot u}) &= \sum_{q=1}^{M_{iju}} c_{tijq} \bar{p}_{tijq} = Y_{tij \cdot u} \\ &= \text{total number of calls in error in} \\ &\quad M_{iju} \text{ entities of type u} \\ \hat{p}_{tij \cdot u} &= \frac{\hat{Y}_{tij \cdot u}}{\hat{X}_{tij \cdot u}} \\ E(\hat{p}_{tij \cdot u}) &= p_{tij \cdot u} = \frac{Y_{tij \cdot u}}{X_{tij \cdot u}} \end{aligned}$$

*It may be noted that in Section 3.1 $\hat{X}_{tij \cdot u}$ were computed on the basis of \bar{n}_{tijq} as the average number of calls per line in (i,j). But, $\hat{X}_{tij \cdot u}$ does not depend on \bar{n}_{tijq} ; it is an estimate of monthly call volume in M_{iju} entities.

A BETTER APPROXIMATION FOR $v(\hat{p}_{ti\cdot\cdot u})$

Recall that

$$\hat{p}_{ti\cdot\cdot u} = \sum_j \left(\frac{\hat{x}_{tij\cdot u}}{\hat{x}_{ti\cdot\cdot u}} \right) \hat{p}_{tij\cdot u} \quad (B1)$$

It is reasonable to assume that $\left(\frac{\hat{x}_{tij\cdot u}}{\hat{x}_{ti\cdot\cdot u}} \right)$ and $\hat{p}_{tij\cdot u}$ are independent. Assuming this independence, it can be shown that

$$\begin{aligned} v(\hat{p}_{ti\cdot\cdot u}) &= \sum_j \left(\frac{\hat{x}_{tij\cdot u}}{\hat{x}_{ti\cdot\cdot u}} \right)^2 \cdot v(\hat{p}_{tij\cdot u}) \\ &+ \sum_j (\hat{p}_{tij\cdot u})^2 \cdot v\left(\frac{\hat{x}_{tij\cdot u}}{\hat{x}_{ti\cdot\cdot u}}\right) \\ &+ \sum_j v\left(\frac{\hat{x}_{tij\cdot u}}{\hat{x}_{ti\cdot\cdot u}}\right) \cdot v(\hat{p}_{tij\cdot u}) \quad (B2) \end{aligned}$$

It can be shown that approximately,

$$v\left(\frac{\hat{x}_{tij\cdot u}}{\hat{x}_{ti\cdot\cdot u}}\right) = \frac{v(\hat{x}_{tij\cdot u})}{\hat{x}_{ti\cdot\cdot u}^2} + \left(\frac{\hat{x}_{tij\cdot u}^2}{\hat{x}_{ti\cdot\cdot u}^4}\right) v(\hat{x}_{ti\cdot\cdot u}) \quad (B3)$$

Since $\hat{x}_{ti\cdot\cdot u} = \sum_j \hat{x}_{tij\cdot u}$, it is sufficient to compute $v(\hat{x}_{tij\cdot u})$. Recall that

$$\hat{x}_{tij\cdot u} = \frac{1}{w_{ij\cdot}} \left[\sum_{q=1}^{m_{ij}} \left(\frac{v_{ij\cdot q}}{v_{ijq}} \right) (c_{tijqu}) (w_{ijq}) \right] \quad (B4)$$

as given in Section 3.2.

$$\therefore v(\hat{x}_{tij\cdot u}) = \frac{1}{(w_{ij\cdot})^2 (w_{ij\cdot} - 1)}$$

$$\cdot \sum_{q=1}^{m_{ij}} \left[\left(\frac{v_{ij\cdot q}}{v_{ijq}} \right) c_{tijqu} - \hat{x}_{tij\cdot u} \right]^2 w_{ijq} \quad (B5)$$

Now putting together all the pieces,

$$\begin{aligned} v(\hat{p}_{ti\cdot\cdot u}) &= \sum_j \left(\frac{\hat{x}_{tij\cdot u}}{\hat{x}_{ti\cdot\cdot u}} \right)^2 v(\hat{p}_{tij\cdot u}) + \sum_j (\hat{p}_{tij\cdot u})^2 \\ &\cdot \left[\frac{v(\hat{x}_{tij\cdot u})}{\hat{x}_{ti\cdot\cdot u}^2} + \frac{\hat{x}_{tij\cdot u}^2}{\hat{x}_{ti\cdot\cdot u}^4} \cdot \sum_j v(\hat{x}_{tij\cdot u}) \right] \\ &+ \sum_j v(\hat{p}_{tij\cdot u}) \cdot \left[\frac{v(\hat{x}_{tij\cdot u})}{\hat{x}_{ti\cdot\cdot u}^2} + \frac{\hat{x}_{tij\cdot u}^2}{\hat{x}_{ti\cdot\cdot u}^4} \sum_j v(\hat{x}_{tij\cdot u}) \right] \quad (B6) \end{aligned}$$

where $v(\hat{x}_{tij\cdot u})$ is given by Equation (B5). To see that the first term on the right hand side of Equation (B6) is the dominant one, examine the following re-expression of (B6):

$$\begin{aligned} v(\hat{p}_{ti\cdot\cdot u}) &= \sum_j \left(\frac{\hat{x}_{tij\cdot u}}{\hat{x}_{ti\cdot\cdot u}} \right)^2 \cdot (\hat{p}_{tij\cdot u})^2 \cdot \frac{v(\hat{p}_{tij\cdot u})}{(\hat{p}_{tij\cdot u})^2} \\ &+ \sum_j \left(\frac{\hat{x}_{tij\cdot u}}{\hat{x}_{ti\cdot\cdot u}} \right)^2 \cdot (\hat{p}_{tij\cdot u})^2 \\ &\cdot \left[\frac{v(\hat{x}_{tij\cdot u})}{(\hat{x}_{tij\cdot u})^2} + \left(\frac{1}{\hat{x}_{ti\cdot\cdot u}} \right)^2 \cdot \sum_j \left(\frac{v(\hat{x}_{tij\cdot u})}{(\hat{x}_{tij\cdot u})^2} \right) \right. \\ &\cdot \left. \left(\hat{x}_{tij\cdot u}^2 \right) \right] \\ &+ \sum_j \left(\frac{\hat{x}_{tij\cdot u}}{\hat{x}_{ti\cdot\cdot u}} \right)^2 (\hat{p}_{tij\cdot u})^2 \left[\frac{v(\hat{p}_{tij\cdot u})}{(\hat{p}_{tij\cdot u})^2} \right] \quad (B7) \\ &\cdot \left[\frac{v(\hat{x}_{tij\cdot u})}{(\hat{x}_{tij\cdot u})^2} + \left(\frac{1}{\hat{x}_{ti\cdot\cdot u}} \right)^2 \cdot \sum_j \left(\frac{v(\hat{x}_{tij\cdot u})}{(\hat{x}_{tij\cdot u})^2} \right) \right. \\ &\cdot \left. \left(\hat{x}_{tij\cdot u}^2 \right) \right] \end{aligned}$$

Since, in general $\hat{p}_{tij\cdot u}$ is much smaller than $\frac{v(\hat{p}_{tij\cdot u})}{(\hat{p}_{tij\cdot u})^2}$, $\frac{\hat{x}_{tij\cdot u}}{\hat{x}_{ti\cdot\cdot u}}$, $\frac{v(\hat{x}_{tij\cdot u})}{(\hat{x}_{tij\cdot u})^2}$ would be much larger than $\frac{v(\hat{x}_{tij\cdot u})}{(\hat{x}_{tij\cdot u})^2}$ and the first term on the left hand side of Equation (B7) is the dominant one.