J.C. Koop and M.P. Singh Dominion Bureau of Statistics Ottawa

# 1. Introduction

Consider a universe where the units are grouped into a number of class intervals or strata according to the variate values of an important characteristic of each of these units. Hereafter the variate value of such a characteristic will be termed the stratification variable. Due to changes in the stratification variable each group of units assigned to a particular stratum (out of a total of L) at a point or epoch of time before a survey is now distributed in a particular manner over L+l strata, hereafter to be termed transitional strata, at the time of the survey, the (L+1)th stratum consisting of units having one or more special characteristics, e.g. going out of business in the context of a farm economy. In this paper we shall define the special characteristic for those units found in each of the (L+1)th strata as one of having zero variate values for all variables of interest to the survey.

At a point or epoch of time before a survey let  $N_h$  be the number of units in the hth stratum (h = 1, 2, ..., L). This is known to the sample surveyor. At the time of the survey the stratification variable of some of the units of the hth stratum changes, and  $N_h$  out of  $N_h$  units are found in the j<sup>th</sup> transitional stratum (j = 1, 2, ..., L+1); the sample surveyor has evidence of this change only from his sample. Let  $x_{hjk}$  be a variate value of interest of the h<sup>th</sup> original stratum (h = 1, 2, ..., L; $j = 1, 2, ..., L+1; k = 1, 2, ..., N_h)$ . The problem is to estimate

$$\begin{array}{c} L \quad L+1 \quad N_{hj} \\ T = \sum \quad \sum \quad \sum \quad x_{hjk}, \quad (1) \\ h=1 \quad j=1 \quad k=1 \quad hjk, \end{array}$$

on the basis of a stratified sample of  $n_1$ ,  $n_2$ ,

..., n<sub>T</sub> units drawn by simple random sampling

without replacement from each respective original stratum.

If 
$$\bar{\bar{x}}_{hj} = \sum_{k=1}^{N_{hj}} x_{hjk} / N_{hj}$$
, for all (h,j),

is the mean of the j<sup>th</sup> transitional stratum of the h<sup>th</sup> original stratum, then (1) may be rewritten as

$$\begin{array}{c} \mathbf{L} \quad \mathbf{L} \\ \mathbf{T} = \sum \sum \sum N_{hj} \, \overline{\mathbf{x}}_{hj}, \qquad (1.1) \\ h=1 \, j=1 \end{array}$$

since  $\bar{x}_{hj} = 0$  for h = 1, 2, ..., L and j = L+1. Further with respect to the original

stratification if 
$$\bar{\bar{x}}_{h} = \sum_{\substack{h \\ j=1}}^{L} N_{hj} \bar{\bar{x}}_{hj} / N_{h}$$
, then

(1.1) may also be rewritten as

$$T = \sum_{h=1}^{L} N_h \overline{\bar{x}}_h.$$
 (1.2)

### 2. Estimation Problems

The stratified sample shows that  $n_{hj}$  out of  $n_h$  units, drawn from the h<sup>th</sup> original stratum, are found in the j<sup>th</sup> transitional stratum in regard to which

$$0 \leq n_{hj} \leq n_{h}$$
, for all  $(h,j)$ . (2)

The universe total T may be estimated

- (i) by standard theory, effectively ignoring the information provided by (2) regarding the transition of units from stratum to stratum, or
- (ii) by taking into account these transitions.

 $\frac{\text{Estimation by standard theory.}}{\text{the usual estimator of T is}}$ 

$$\hat{\mathbf{T}} = \sum_{h=1}^{L} N_{h} \bar{\mathbf{x}}_{h}, \qquad (3)$$

where  $\bar{\mathbf{x}}_{h}$  is the mean of the sample units in stratum h. Alternatively in (1.1) if unbiased estimates of  $N_{hj}$  and  $\bar{\mathbf{x}}_{hj}$  are substituted, viz.

$$\hat{N}_{hj} = \frac{n_{hj}}{n_h} N_h, \qquad (4)$$

and

$$\bar{\bar{x}}_{hj} = \sum_{k=1}^{n_{hj}} x_{hjk} / n_{hj}, \qquad (5)$$

for all relevant (h,j), then it is not difficult to see that the resulting estimator does not involve the  $n_{hj}$ 's and is identical to (3). The variance of this estimator is

$$V(\hat{T}) = \sum_{h=1}^{L} N_{h}^{2} \frac{S_{h}^{2}}{n_{h}} (1 - \frac{n_{h}}{N_{h}}), \qquad (6)$$

where

$$S_{h}^{2} = \sum_{j=1}^{N} \sum_{k=1}^{N} (x_{hjk} - \overline{x}_{h})^{2} / (N_{h} - 1) \text{ for all } h.$$

Estimation taking into account transition probabilities. The probability of transition of a unit in stratum h to the jth transitional stratum is  $p_{hj} = N_{hj}/N_h$ , and on the basis of sample data an unbiased estimate of this transition probability is  $p_{hj} = n_{hj}/n_h$  for all (h,j). Thus we have a matrix, with L+1 rows and L columns, of the estimated transition probabilities

$$(\hat{p}_{hj}).$$
 (7)

This matrix multiplied by the column vector of the number of units in the L strata at the point of epoch of time before the survey, i.e.,  $\{N_1, N_2, \ldots, N_L\}$ , yields a vector of L+1 unbiased estimates of the number of units in the L+1 strata at the time of the survey, the number in the j<sup>th</sup> stratum being

$$\hat{\mathbf{N}}_{,j} = \sum_{\mathbf{h}=1}^{L} \hat{\mathbf{p}}_{\mathbf{h}j} \mathbf{N}_{\mathbf{h}}, \qquad (8)$$

which is simply the sum over h of the estimate given by (4). By classical theory the variance of  $\hat{N}_{i}$ , for all j, is given by

$$V(\hat{N}_{.j}) = \sum_{h=1}^{L} N_{h}^{2} \frac{p_{hj}(1-p_{hj})}{n_{h}} (\frac{N_{h}-n_{h}}{N_{h}-1}). \quad (9.1)$$

Also

$$\operatorname{Cov}(\hat{\mathbb{N}}_{j}, \hat{\mathbb{N}}_{j}, ) = -\sum_{h} N_{h}^{2} p_{hj} p_{hj}, (\frac{N_{h}-n_{h}}{N_{h}-1})/n_{h}$$
  
for all  $j \neq j'$ . (9.2)

On the basis of the estimated stratum numbers given by (8), a biased but consistent estimator of T is

$$\mathbf{T}' = \sum_{j=1}^{L} \hat{\mathbf{N}} \cdot \mathbf{j} \left( \sum_{h=1}^{L} \mathbf{n}_{hj} \cdot \overline{\mathbf{x}}_{hj} / \sum_{h=1}^{L} \mathbf{n}_{hj} \right). \quad (10)$$

The loss in efficiency of the original stratification is to some extent regained by restratification through pooling each set of L transitional strata relating to the same range of x-values; this is the justification for constructing an alternative estimator T' given by (10).

We shall derive briefly an expression for the bias of T'. Substituting the expression for  $\hat{N}_{,j}$  in (10), and denoting the random variables specified in (2) by  $n_{,j}$  we find

$$E(\mathbf{T}') = E\{E(\mathbf{T}' | \underline{p}_{hj})\}$$
$$= E\{\sum_{j=1}^{L} (\sum_{h=1}^{L} \frac{n_h}{n_h} n_{hj}) (\frac{\sum_{j=1}^{\Sigma} n_{hj} \overline{x}_{hj}}{\sum_{h=1}^{\Sigma} n_{hj}})\}$$

$$= \sum_{j=1}^{L} \sum_{h=1}^{L} \frac{N_{h} \overline{x}_{hj}}{n_{h}} E(\frac{n_{hj}^{2}}{L})$$

$$+ \sum_{j=1}^{L} \sum_{h\neq h'}^{L} \frac{N_{h} \overline{x}_{h'j}}{n_{h}} E(\frac{n_{hj}n_{h'j}}{L}). \qquad (11)$$

In order to evaluate the expectations of ratios of functions of random variables found in (11) we need the following result due to one of us. If X and Y are pairs of random variables with X assuming all values but not zero, then

$$E(\frac{Y}{X}) = \frac{E(Y)}{E(X)} - \frac{Cov(Y,X)}{E^{2}(X)} + \frac{1}{E^{2}(X)} E\{\frac{Y}{X}(X-E(X))^{2}\}.$$
 (12)

With this result, we find

$$E(\frac{n_{hj}^{2}}{\sum n_{hj}}) = \frac{E(n_{hj}^{2})}{E(\sum n_{hj})} - \frac{Cov(n_{hj}^{2}, \sum n_{hj})}{E^{2}(\sum n_{hj})} + \frac{1}{E^{2}(\sum n_{hj})} E\{\frac{n_{hj}^{2}}{\sum n_{hj}} + \frac{1}{E^{2}(\sum n_{hj})} E\{\frac{n_{hj}^{2}}{\sum n_{hj}} + \frac{1}{E^{2}(\sum n_{hj})} + \frac{1}{E^{2}($$

and

4

$$E\left(\frac{n_{hj}n_{j}'}{\sum n_{hj}}\right) = \frac{E(n_{hj}n_{j}')}{E(\sum n_{hj})} - \frac{Cov(n_{hj}n_{j}')}{E^{2}(\sum n_{hj})} + \frac{1}{E^{2}(\sum n_{hj})} + \frac{1}{E^{2}(\sum n_{hj})} E\left(\frac{n_{hj}n_{j}'}{\sum n_{hj}}\right) = E\left(\frac{n_{hj}n_{j}'}{\sum n_{hj}}\right) + \frac{1}{E^{2}(\sum n_{hj})} + \frac{1}{E^{2}(\sum n_{hj})} + \frac{1}{E^{2}(\sum n_{hj})} + \frac{1}{E^{2}(\sum n_{hj})} = E\left(\frac{n_{hj}n_{j}'}{\sum n_{hj}}\right) + \frac{1}{E^{2}(\sum n_{hj})} + \frac{1}{E^{2}(\sum n_{$$

Further details of the derivation of E(T') are found in the appendix. Considering only the leading terms of (13) and (14) we find

$$E(T') = T + \frac{1}{n} \left[ \sum_{j=1}^{L} \{ \sum_{h>h}, N_{h}N_{h}, P_{hj}P_{h'j} \right]$$
$$(\bar{x}_{hj} - \bar{x}_{h'j}) \left( \frac{n_{h}}{N_{h}} - \frac{n_{h'}}{N_{h'}} \right) / \left( \sum_{h=1}^{L} \alpha_{h}P_{hj} \right)$$
$$- \sum_{h=1}^{L} N_{hj} \bar{x}_{hj} (1 - P_{hj}) \left( \frac{N_{h} - n_{h}}{N_{h} - 1} \right) / \left( \sum_{h=1}^{L} \alpha_{h}P_{hj} \right) \}, (15)$$

where  $\alpha_h = n_h/n$  is the proportion of the sample allocated to stratum h. Now  $\bar{x}_{hj}$  is the mean of all the  $N_{hj}$  units of the h<sup>th</sup> original stratum which have moved to the j<sup>th</sup> transitional stratum. Similarly  $\bar{x}_{h'j}$  is the mean of all the  $N_{h'j}$  units of the h'<sup>th</sup> original stratum which have moved to the j<sup>th</sup> transitional stratum. Since by this procedure of restratification the range of x-values in each of the j<sup>th</sup> transitional strata is controlled, the differences  $\bar{x}_{hj} = \bar{x}_{h'j}$ , for all h>h', could not be considerable. Further it is unlikely that the signs of such differences are the same unless the class interval of the stratification variable is wide. Indeed (15) suggests that the bias, E(T')-T, can be reduced by narrowing the width of the transitional strata so as to make the  $\bar{x}$ 's for each j as equal as possible, in which situation this bias will be negligible when n is large.

The first steps in the derivation of the variance of T' are as follows:

$$V(\mathbf{T}') = V\{E(\mathbf{T}'|\underline{n}_{hj})\} + E V(\mathbf{T}'|\underline{n}_{hj})\}$$
  
=  $V\{\sum_{j=1}^{L} (\sum_{h=1}^{L} \frac{N_h}{n_h} n_{hj}) (\frac{\sum_{j=1}^{L} n_{hj} \overline{x}_{hj}}{\sum_{h=1}^{L} n_{hj}})\}$ 

+ 
$$E\left\{\sum_{j=1}^{L}\left(\sum_{h=1}^{L}\frac{N_{h}}{n_{h}}n_{hj}\right)^{2}\frac{1}{(\sum n_{hj})^{2}}\cdot\right.$$
  
 $\left.\sum_{h=1}^{L}n_{hj}^{2}\frac{S_{hj}^{2}}{n_{hj}}\left(1-\frac{n_{hj}}{N_{hj}}\right)\right\},$  (16)

where  $S_{hj}^2 = \sum_{k=1}^{N_{hj}} (x_{hjk} - \overline{\overline{x}}_{hj})^2 / (N_{hj} - 1)$  for all

(h,j). The first part of the variance function V{ } may be ascribed to movement of some of the units of the original strata to their respective transitional strata resulting in the  $\bar{\bar{x}}_{h,j}$ 's. This variance function involves the products of the estimates of the transition probabilities  $n_{hj}/n_h$  and the corresponding transitional stratum weights  $n_{hj} / \sum_{h} n_{hj}$ , and also involves the relevant cross-product terms. It can be shown that the variances and covariances of all these product and cross-product terms are of order  $1/n_h^2$  for all h so that there is control on this part of the variance through restratification. The second part of the variance function is due to the variation of the x's in the transitional strata. We again remark that this source of variation can be controlled by narrowing the width of the transitional strata.

The variance of the unbiased estimator given by (6) can be rewritten as

$$\begin{aligned} \mathbf{v}(\hat{\mathbf{T}}) &= \sum_{h=1}^{L} N_{h}^{2} \left(1 - \frac{n_{h}}{N_{h}}\right) \frac{1}{n_{h}} \sum_{j=1}^{L} \left(\frac{N_{hj}}{N_{h}-1}\right) \left(\overline{\mathbf{x}}_{hj} - \overline{\mathbf{x}}_{h}\right)^{2} \\ &+ \sum_{h=1}^{L} \frac{N_{h}^{2}}{n_{h}} \left(1 - \frac{n_{h}}{N_{h}}\right) \left(\frac{N_{h,1}+1}{N_{h}-1}\right) \overline{\mathbf{x}}_{h}^{2} \\ &+ \sum_{h=1}^{L} N_{h}^{2} \left(1 - \frac{n_{h}}{N_{h}}\right) \frac{1}{n_{h}} \sum_{j=1}^{L} \left(\frac{N_{hj}-1}{N_{h}}\right) S_{hj}^{2} \cdot (17) \end{aligned}$$

The expression speaks for itself. There are reasons to believe that  $V(\hat{T})$  can be greater than V(T'), because of the middle term in (17).

The variance expression for T' is developed to this point just to gain an insight into the features commented upon in the foregoing account. An exact expression for it will involve all the transition probabilities, and is not useful for estimation purposes because of its complexity.

The problem of variance estimation can be resolved by redesigning the survey so as to have two or more independent stratified samples, each yielding an independent estimate of T. When the total sample consists of two independent replicates, then the variance of the mean  $\overline{T}=(T_1'+T_2')/2$  is given by

$$v(\bar{T}) = (T_1' - T_2')^2/4,$$
 (18)

despite the complexity of (16).

### 3. Remarks

The problem has been considered in the context of one-stage stratified sampling partly because it is simple, the strata corresponding to the "states" in the theory of Markov chains, and partly because the problem was originally seen in relation to this sampling design. The problem of estimation can be considered almost along the same lines for more ramified sampling designs.

<u>Acknowledgement</u>. We are grateful to Dr. I.P. Fellegi for some critical remarks which have improved the presentation of this paper. To our colleagues in the Agriculture Division of this Bureau, we are also grateful for many enlightening discussions on the practical problems of estimation in livestock and poultry sample surveys, which have provided the impulse for the estimation problem considered in this paper.

#### Reference

Koop, J.C. On the derivation of expected value and variance of ratios without the use of infinite series expansions. To appear in Metrika.

## Appendix

<u>Derivation of E(T')</u>. For evaluation of terms involved in (13) and (14) we note the following.

$$\begin{split} \mathbb{E}(\mathbf{n}_{hj}) &= \mathbf{n}_{h}\mathbf{p}_{hj}, \quad \mathbb{E}(\mathbf{n}_{hj}^{2}) = \mathbf{n}_{h}^{2}\mathbf{p}_{hj}^{2} + \mathbf{n}_{h}\mathbf{p}_{hj}\mathbf{q}_{hj}\mathbf{f}_{h}, \\ \mathbb{E}(\mathbf{n}_{hj}\mathbf{n}_{h'j}) &= \mathbf{n}_{h}\mathbf{n}_{h'}\mathbf{p}_{hj}\mathbf{p}_{h'j}, \\ & \text{where } \mathbf{q}_{hj} = \mathbf{1} - \mathbf{p}_{hj}, \quad \mathbf{f}_{h} = (\mathbf{1} - \frac{\mathbf{n}_{h} - \mathbf{1}}{\mathbf{N}_{h} - \mathbf{1}}), \\ \mathbb{Cov}(\mathbf{n}_{hj}^{2}, \sum_{h=1}^{L} \mathbf{n}_{hj}) = \mathbb{E}(\mathbf{n}_{hj}^{3}) - \mathbb{E}(\mathbf{n}_{hj}^{2}) \mathbb{E}(\mathbf{n}_{hj}) \\ & \stackrel{\cdot}{\neq} \mathbf{n}_{h}^{2}\mathbf{p}_{hj}^{2} (3 - \mathbf{q}_{hj}\mathbf{f}_{j}) - \mathbf{n}_{h}\mathbf{p}_{hj} (2 - 3\mathbf{q}_{hj}\mathbf{f}_{h}), \end{split}$$

and  

$$Cov(n_{hj}n_{h'j}, \sum_{h=1}^{L} n_{hj}) = Cov(n_{hj}n_{h'j}, n_{hj}) + Cov(n_{hj}n_{h'j}, n_{h'j}) + Cov(n_{hj}n_{h'j}, n_{h'j}) + Cov(n_{hj}n_{h'j}, n_{h'j}) + Cov(n_{hj}n_{h'j}, n_{h'j}) + Cov(n_{hj}) + Cov(n_{hj}) + Cov(n_{hj}) + Cov(n_{h'j}) + Cov(n_{$$

$$-\frac{1}{n^{2}}\sum_{j=1}^{L}\sum_{h\neq h}^{L},\frac{N_{h}\overline{x}_{h}j}{(\sum \alpha_{h}p_{hj})^{2}},$$
$$\{n_{h}p_{h}j_{h}p_{hj}(f_{h}q_{hj}+f_{h}q_{h}j)-\frac{R_{hj}^{(2)}}{n_{h}}\},$$

where  $\alpha_h = n_h/n$  is the proportion of sample allocated to stratum h. Ignoring the contributions of third and fourth term for large sample size and adding and subtracting

$$\begin{array}{c} {}^{\mathrm{L}}_{\Sigma} ( \mathbb{N}_{\mathrm{h}} \overline{\overline{\mathbf{x}}}_{\mathrm{h}j} \mathbf{p}_{\mathrm{h}j} ) ( \mathbb{n}_{\mathrm{h}}, \mathbb{p}_{\mathrm{h}}, \mathbf{j} ) / \sum_{\mathrm{h}=1}^{\mathrm{L}} \mathbb{n}_{\mathrm{h}} \mathbb{p}_{\mathrm{h}j} \\ \mathbb{h}_{\mathrm{h}}^{\frac{1}{2}} \mathbb{h}^{\mathrm{h}} \mathbf{h}^{\frac{1}{2}} ( \mathbb{h}_{\mathrm{h}}, \mathbb{h}_{\mathrm{h}}, \mathbb{h}_{\mathrm{h}}) \\ \end{array}$$

to the first term we can write

$$E(T') = \sum_{j=1}^{L} \sum_{h=1}^{L} N_{hj} \overline{\bar{x}}_{hj}$$

$$+ \frac{1}{n} \sum_{j=1}^{L} \{\sum_{h=h}^{L} N_{hj} (\overline{\bar{x}}_{h'j} - \overline{\bar{x}}_{hj}) (\frac{n_{h'}p_{h'j}}{L})$$

$$+ \sum_{h=1}^{L} N_{hj} \overline{\bar{x}}_{hj} (\frac{q_{hj}f_{j}}{L}),$$

$$\sum_{h=1}^{L} \alpha_{h}p_{hj}$$

which can also be expressed in the form given in (15).

.