Gad Nathan, Hebrew University, Jerusalem and Central Bureau of Statistics. Israel*

0. Introduction and Summary: The basic response error model developed by the U.S. Bureau of the Census [5] has been generalized to the multivariate case by Koch [4]. The possibilities of application of the model to complex sample designs and to complex estimators are, however, limited by the requirement that all components of the vector random variable are measured for the same sample. This is not necessarily the case for many applications, such as difference, ratio or regression estimation, or panel-type sample designs, where component estimators may be based on different samples (e.g. sub-samples, partially matched samples, etc.). In the following, the basic multivariate response error model is extended to cover the case where each component of the vector variable is measured in a possibly different sample, all selected from the same finite population. The extended model is then applied to difference and ratio estimation in different response error structure situations and to the case of sampling on two occasions.

1. The Model: The formulation of the model follows that of Koch [4]. For each unit in the population, i=1,...,N , let there be defined a p-component vector random variable, X1+, by:

 $y_{it}^{\prime} = (y_{it}^{(1)}, y_{it}^{(2)}, \dots, y_{it}^{(p)})$, where t indexes the sequence of repeated trials.

Any given sample is defined by the indicator random variable:

$$U_{i}^{(j)} = \begin{cases} 1 : \text{ if unit } i \text{ is in the sample} \\ \text{for the } j\text{-th component} \\ 0 : \text{ otherwise.} \end{cases}$$

So as to simplify the presentation, simple random sampling without replacement will be assumed for each of the component samples. Thus, if n,

is the sample size for component j, then:

 $E\{U_{i}^{(j)}\}=n_{j}/N \text{ and: } E\{U_{i}^{(j)}U_{i'}^{(j)}\}=n_{j}(n_{j}-1)/[N(N-1)],$ $(i\neq i') \text{ The relationship between the}$ samples is defined by: $P\{U_{i}^{(j)}=U_{i}^{(j')}=1\}=E\{U_{i}^{(j)}U_{i}^{(j')}\}=\begin{cases}v_{jj}, : i=i'\\v_{ji}, : i\neq i', \\v_{i+1}, : i\neq i', \end{cases}$

where it is assumed that the expectation depends only on whether i=i' and not on the specific values of i and i'. If we define: $n_{jj}^{n_j}, j_{j'}^{n_j}, j_{j'}^{n_j}$, (where $n_{jj}^{n_j} = n_j$ and it is assumed for the time being that $v_{jj}, \neq 0$), then

it is easy to see that:

$$w_{jj'} = \frac{n_j n_{j'} (n_{jj'} - 1)}{n_{jj'} N(N-1)}$$
.

The statistic considered will be the sample mean: $\overline{y}_{t}^{\prime} = (\overline{y}_{t}^{(1)}, \overline{y}_{t}^{(2)}, \dots, \overline{y}_{t}^{(p)}),$ where:

 $\begin{array}{c} \overline{y}_{t}^{(j)} \sum U_{i}^{(j)} Y_{it}^{(j)} / n_{j} , \text{ as an estimate of the popula-} \\ tion mean: \quad \overline{\chi}' = (1/N) \sum \chi_{i}' , \text{ where: } \chi_{i}' = \mathbb{E}_{t} \{\chi_{it}'\} \\ i = 1 \end{array}$

is the expected response for the 1-th element over all trials. It will be assumed that there is no response bias, so that:

 $\begin{array}{l} Y_{i}^{(j)} = E_{t} \{Y_{it}^{(j)}\} = E_{t} \{Y_{it}^{(j)} | U_{i}^{(j)} = 1\} \\ \text{covariance matrix of } \overline{y}_{t} \\ \end{array}$ can be decomposed, as usual, as follows:

 $\sum_{j=1}^{n} \frac{(j)}{1} \frac{y^{(j)}}{n}$ be the sample mean of the expec-

ted values for the j-th component, this decomposition reflects the three sources of variation as follows:

Response Variance: $RY = E \{ (\overline{y}, -\overline{y}) (\overline{y}, -\overline{y})' \}$

Sampling Variance: $SY=E_{+}\{(\overline{y}-\overline{y})(\overline{y}-\overline{y})'\}$

Interaction:

$$2\mathcal{U} = \mathcal{U} + \mathcal{U} +$$

In the following, $v^{(j,j')}$, the (j,j') element of χ , will be expressed by further decomposing each of its components:

$$V^{(j,j')}_{RV}^{(j,j')}_{+SV}^{(j,j')}_{+2} \overline{IV}^{(j,j')}$$
. (2)

Define the simple response variance as:

$$SRV_{i}^{(jj')} = E_{t} \{ (Y_{it}^{(j)} - Y_{i}^{(j)}) (Y_{it}^{(j')} - Y_{i}^{(j')}) | U_{i}^{(j)} = U_{i}^{(j')} = 1 \}$$

and
SRV_{ij'}^{(jj')} = (1/N) $\Sigma_{i} SRV_{i}^{(jj')}$ (3)
(3)

and, similarly the correlated response variance as:

$$CRV_{ii}^{(jj')} = E_{t} \{ (Y_{it}^{(j)} - Y_{i}^{(j)}) (Y_{i't}^{(j')} - Y_{i'}^{(j')}) U_{i}^{(j)} U_{i'}^{(j')} \}$$

$$and:$$

$$CRV_{ij'}^{(jj')} = \frac{1}{N(N-1)} \sum_{i \neq i'} CRV_{ii'}^{(jj')} ,$$

$$(4)$$

^{*}This research was carried out while the author was Visiting Associate Professor at the University of North Carolina at Chapel Hill and at the Research Triangle Institute and was sponsored by U.S. Bureau of the Census Contract No. 1-35096.

(where the conditional expectations are well defined, since we assumed

 v_{jj} ,=P{ $U_{i}^{(j)}=U_{i}^{(j')}=1$ } ≠0 and it can easily be seen that w_{jj} ,=P{ $U_{i}^{(j)}=U_{i'}^{(j')}=1$ } ≠0, except in the trivial case $n_{j}=n_{j}$,=1).

The response variance component can be shown to be:

If we define the simple sampling variance as:

$$ssv^{(jj')} = \frac{1}{N-1} \Sigma_{i}(Y_{i}^{(j)} - \overline{Y}^{(j)})(Y_{i}^{(j')} - \overline{Y}^{(j')}), \quad (6)$$

the sampling variance component can be written as

$$sv^{(jj')} = (1/n_{jj'})(1 - \frac{n_{jj'}}{N})ssv^{(jj')}.$$
 (7)

The <u>interaction component</u>, IV^(jj'), which reflects the inter-relationship of sampling and response errors, is non-zero if:

$$Y_{ii'}^{(jj')} = E_{t} \{Y_{it}^{(j)} | U_{i}^{(j)} = U_{i}^{(j')} = 1\} \neq Y_{i}^{(j)}$$

If we define simple interaction between response and sampling deviations for the same units as:

$$siv^{(jj')} = (1/N) \Sigma_{i} (Y_{ii}^{(jj')} - Y_{i}^{(j)}) Y_{i}^{(j')}$$

and: $\overline{siv}^{(jj')} = 1/2 (siv^{(jj')} + siv^{(j'j)}) ; (8)$

and define the simple correlated interaction between response and sampling deviations for different units as:

$$sciv^{(jj')} = \frac{1}{N(N-1)} \sum_{i \neq i} (Y_{ii}^{(jj')} - Y_{i}^{(j)}) Y_{i}^{(j')}$$

and: $\overline{sciv}^{(jj')} = \frac{1}{2} (sciv^{(jj')} + sciv^{(j'j)}) . \quad (9)$
Then: $\overline{iv}^{(jj')} = \frac{1}{n_{jj}} (\overline{siv}^{(jj')} + n_{ij}^{(jj')}) . \quad (10)$

Substituting (5), (7) and (10) in (2) we obtain:

$$v^{(jj')} = \frac{1}{n_{jj'}} (SRV^{(jj')} + (n_{jj'} - 1)CRV^{(jj')}) + (1 - n_{jj'} / N)SSV^{(jj')} + 2[\overline{SIV}^{(jj')}) + (n_{jj'} - 1)\overline{SCIV}^{(jj')}] , \qquad (11)$$

where the last two terms drop out if

 $Y_{11}^{(jj')}=Y_{1}^{(j)}$ for all (i,i'). It should be noted that this expression has the same form as that given by Koch [4] for the case where all components are measured on the same sample, with only the sample size n replaced by n_{11} . Thus, methods proposed by Chai [2] and jj'Bailar and Dalenius [1] to estimate the population parameters on the basis of a single sample can be applied.

For the special case, $v_{ij} = 0$ ($j \neq j'$), i.e.

non-overlapping samples for the j-th and j'-th component, the resulting modification of the decomposition, (for $j \neq j'$), is:

$$v^{(jj')}=SCRV^{(jj')}+IRV^{(jj')}-\frac{1}{N}SSV^{(jj')}+2SCIV^{(jj')},$$
(12)

which is the limit of (11) as $1/n_{jj}$, goes to zero. It should be noted that $v^{(jj')}$ is independent of sample size, in this case.

2. <u>Application to Complex Estimators</u>: The above model can easily be applied to a variety of sampling designs with respect to the relationship between samples for different components. In the following, the application of the extended model to complex estimators is considered. For the sake of algebraic simplicity, we shall assume no interaction between sampling and response deviations, in the sense that $Y_{11}^{(1)} = Y_{11}^{(1)}$ for all

i,i'=1,...,N and for all j,j'=1,...,p. The results can easily be extended to the case of non-zero interactions. The variance of a linear combination of the sample means: - -(i) - -(i)

 $\overline{y}_{t\xi} = \sum_{j} \ell_{j} \overline{y}_{t}^{(j)}$, as an estimate of $\overline{Y}_{\xi} = \sum_{j} \ell_{j} \overline{Y}^{(j)}$, will be:

$$\operatorname{var}(\overline{y}_{t_{k}})=\Sigma_{j,j}, \mathfrak{l}_{j}\mathfrak{l}_{j}, \nabla^{(jj')} \qquad (13)$$

Similarly the variance of an analytical function of the sample means $g(\overline{y}_t^{(1)}, \dots, \overline{y}_t^{(p)})$ can be approximated by the appropriate Taylor expansion.

We shall consider two variables, X and Y, with observable values at the t-th trial for the i-th unit X_{it} , Y_{it} . Define

 $X_i = E_t(X_{it}); Y_i = E_t(Y_{it})$ and let $\overline{X} = (1/N)\Sigma_i X_i$, $\overline{Y} = (1/N)\Sigma_i Y_i$, $\overline{X}_t = (1/N)\Sigma_i X_{it}$ and $\overline{Y}_t = (1/N)\Sigma_i Y_{it}$, be the population means of the expected values and of the values for the t-th trial respectively. Let:

$$V_{xy} = (1/N) \sum_{i} E_{t} (X_{it} - X_{i}) Y_{it} - Y_{i}; \text{ and similarly}$$

$$V_{xx} \quad \text{and} \quad V_{yy};$$

$$C_{xy} = \frac{1}{N(N-1)} \sum_{i \neq i} (X_{it} - X_{i}) (Y_{i't} - Y_{i'});$$
and similarly C_{xx} and $C_{yy};$

$$S_{xy} = \frac{1}{N-1} \Sigma_i (X_i - \overline{X}) (Y_i - \overline{Y})$$
; and similarly S_{xx}

values, respectively.

Consider a single simple random sample without replacement of size n, selected from the population of N elements and defined by the indicator random variables U_i (i=1,2,...,N). Let $\overline{x}_t = (1/n)\Sigma_i U_i X_{it}$; $\overline{y}_t = (1/n)\Sigma_i U_i Y_{it}$; $\overline{x} = (1/n)\Sigma_i U_i X_i$; and $\overline{y} = (1/n)\Sigma_i U_i Y_i$ be the sample means for the t-th trial and for the expected In the following Y will be the variable for which an estimate of the population mean of expected values, \overline{Y} , is required. Measurements of Y are available only for the sample elements at a given trial, t, so that only \overline{y}_t is

known. X will be an auxiliary variable for which measurements are available for the whole population, in one of the following alternative ways:

- (a) X is measured at the t-th trial for the population and for the sample. In this case \overline{x}_t and \overline{X}_t are known but not \overline{x} or \overline{X} . This would be the most usual case in practice, giving rise to the difference estimate: $\overline{y}_{Dt1} = \overline{y}_t + k(\overline{x}_t \overline{X}_t)$; or to the ratio estimate: $\overline{y}_{Rt1} = (\overline{y}_t / \overline{x}_t)\overline{X}_t$.
- (b) X is measured at the t-th trial for the sample, so that \overline{x}_t is known, and the population mean of the expected values, \overline{X} , is known (or alternatively \overline{X}_t has no response error), but \overline{x} , the sample mean of expected values, is not known. This case would be rather unusual in practice but may arise if errorless measurements are available for the whole population (or only for its mean) but there is a practical difficulty in matching the sample back, to obtain the values of X_i for the sample

elements. The difference estimate for this case is: $y_{Dt2} = y_t + k(\bar{x}_t - \bar{X})$; and the ratio estimate is: $\bar{y}_{Rt2} = (\bar{y}_y / \bar{x}_t) \bar{X}$.

(c) The expected values of X are known for the whole population and for the sample (or X is assumed to be measured without response error), so that \overline{x} and \overline{X} are known. In this case the difference estimate is $\overline{y}_{Dt3} = \overline{y}_t + k(\overline{x} - \overline{X})$; and the ratio estimate is: $\overline{y}_{R+3} = (\overline{y}_t / \overline{x})\overline{X}$.

The variances of all these estimates can be obtained from (11) and (13), by considering the vector variables: $\chi'_{it} = (Y_{it}, X_{it}, X_{it}), \chi'_{it} =$ (Y_{it}, X_{it}, X_i) or $\chi'_{it} = (Y_{it}, X_i, X_i)$, for (a), (b) or (c), respectively, with sample sizes: $n_1 = n_2 = n$ and $n_3 = N$ and values of n_{jj} . $n_{11} = n_{12} = n_{22} = n$ and $n_{13} = n_{23} = n_{33} = N$.

The variance of the difference estimates is obtained by applying (13) with l=(1,k,-k). The common component of the variances of the three estimates which is independent of k is the variance of the sample mean:

$$var(\overline{y}_{t}) = V_{yy} + (n-1)C_{yy} + (1-n/N)S_{yy}$$
. (14)

If we define the intra-trial correlation as:

$$\delta^{(jj')} = \frac{CRV^{(jj')}}{SRV^{(jj')}} \quad (if \quad SSRV^{(jj')} \neq 0) ; so that:$$

$$\delta_{xx} = C_{xx} / V_{xx} ; \delta_{xy} = C_{xy} / V_{xy} ; \delta_{yy} = C_{yy} / V_{yy} , then$$

minimal variances of the difference estimates are:

$$V_{1} = \min_{k} [n \ var(\overline{y}_{Dt1})] = n \ var(\overline{y}_{t}) - (1-f) \frac{[V_{xy}(1-\delta_{xy})+S_{xy}]^{2}}{V_{xx}(1-\delta_{xx})+S_{xx}}$$
(15)

$$V_{2}^{=\min_{k} [n \ var(\overline{y}_{Dt2})]=n \ var(\overline{y}_{t})} - \frac{[V_{xy} \{1+(n-1)\delta_{xy}\}+(1-f)S_{xy}]^{2}}{V_{xx} \{1+(n-1)\delta_{xx}\}+(1-f)S_{xx}}$$
(16)

$$V_3 = \min_k [n var(\overline{y}_{Dt3})] = n var(\overline{y}_t) - (1 - f)S_{xy}^2 / S_{xx}, (17)$$

where $f = n/N$.

The last terms of the right hand sides of (15)-(17) represent the possible gains over the sample mean from the use of the three difference estimates. If we assume V_{xy} , S_{xy} , $\delta_{xy} > 0$, then for n sufficiently large we will have: $V_2 < V_1$ and $V_2 < V_3$, so that the greatest gain could be made by using \overline{y}_{Dt2} , (i.e. using the trial sample mean, \overline{x}_t , and the population expected mean \overline{X}), even if errorless (i.e. expected) values are available both for the sample and for the population.

If, however, the response correlations, V_{xy} and C_{xy} , are small, relative to the response variances (V_{xx} and C_{xx}), and n is not too large then y_{Dt3} may well have the smallest variance of the three estimates.

The variance of the ratio estimates, \overline{y}_{Rt} , will be approximated by applying (13) with l=(1,-R,R). Thus (approximately): n var (\overline{y}_{Rt}) = n var (\overline{y}_{Rt})

$$\frac{v_{Rt1}^{y} - v_{Rt1}^{y}}{+(1-f) \{R^{2}[v_{xx}^{x}(1-\delta_{xx}^{x})+S_{xx}^{x}] -2R[v_{xy}^{x}(1-\delta_{xy}^{x})+S_{xy}^{x}]\}}$$
(18)

n var
$$(\overline{y}_{Rt2})$$
=n var (\overline{y}_{t})
+R²{ v_{xx} [1+(n-1) δ_{xx}]+(1-f)S_{xx}}
-2R{ v_{xy} [1+(n-1) δ_{xy}]+(1-f)S_{xy}}; (19)

n var
$$(\overline{y}_{Rt3})$$
=n var (\overline{y}_t) +(1-f)[R²S_{xx}-2RS_{xy}] . (20)

The conditions under which the ratio estimates are better than the sample mean, i.e. $var(y_{R+}) < var(y_{+})$, are then:

$$\begin{array}{l} R<2 \; \frac{V_{xy} \left(1-\delta_{xy}\right)+S_{xy}}{V_{xx} \left(1-\delta_{xx}\right)+S_{xx}} \; , \; \text{for } \; \overline{y}_{Rt1} \; ; \\ R<2 \; \frac{V_{xy} \left[1+(n-1)\delta_{xy}\right]+(1-f)S_{xy}}{V_{xx} \left[1+(n-1)\delta_{xx}\right]+(1-f)S_{xx}} \; , \; \text{for } \; \overline{y}_{Rt2} \; ; \end{array}$$

$$R<2 \frac{S_{xy}}{S_{xx}}$$
, for \overline{y}_{Rt3} (21)

If response correlations are small, relative to the correlations between expected values, the conditions (21) for \overline{y}_{Rt1} and \overline{y}_{Rt2} may be more stringent than the condition for \overline{y}_{R+3} and in particular if:

$$\frac{V_{xy}(1-\delta_{xy})}{V_{xx}(1-\delta_{xx})} < \frac{S_{xy}}{S_{xx}} , \text{ for } \overline{y}_{Rt2} \text{ and } \text{if:}$$

$$\frac{V_{xy}[1+(n-1)\delta_{xy}]}{V_{xx}[1+(n-1)\delta_{xx}]} < \frac{S_{xy}}{S_{xx}} , \text{ for } \overline{y}_{Rt1} . (22)$$

The last relationship would always hold if δ <δ for large enough n.

Comparing
$$\overline{y}_{Rt1}$$
 and \overline{y}_{Rt2} , we have
 $var(\overline{y}_{Rt2}) > var(\overline{y}_{Rt1})$ if:
 $R > 2 \frac{V_{xy}[1+(N-1)\delta_{xy}]}{V_{xx}[1+(N-1)\delta_{xx}]} = 2 \frac{E_t\{(\overline{Y}_t - \overline{Y})\overline{K}_t - \overline{X})\}}{E_t\{(\overline{X}_t - \overline{X})^2\}}$. (23)

Thus, if the correlation between the trial means, \underline{X}_{t} and \underline{Y}_{t} , is small, the ratio estimator \overline{y}_{Rt2} would be preferred to \overline{y}_{Rt1} , i.e. the trial population, \overline{X}_{+} , should be used rather than the errorless expected population mean, \overline{X} , for blowing-up the trial sample ratio, y_t/\overline{x}_t , even if \overline{X} is known.

3. Application to Sampling on Two Occasions: Let samples of the same sizes, n , be selected on each of two occasions such that m(<n) elements are matched and measured on both occasions and u(=n-m) are unmatched. Let Y be the variable for the second (current) occasion and X for the first and consider the unbiased estimate from the t-th trial for \overline{Y} :

$$\overline{y}_{t}^{\dagger}=a\overline{y}_{mt}^{\dagger}+(1-a)\overline{y}_{ut}^{\dagger}+b\overline{x}_{mt}^{}-b\overline{x}_{ut}$$
, (24)

where \overline{y}_{mt} , \overline{x}_{mt} are the sample means for the matched part, \overline{y}_{ut} , \overline{x}_{ut} are the sample means for the unmatched parts and a, b are any constants. Set: $\chi'_{it} = (Y_{it}, Y_{it}, X_{it}, X_{it})$ with: $n_1 = n_3 = m$; $n_2 = n_4 = u$; $v_{12} = v_{14} = v_{23} = v_{24} = v_{34} = 0$; $v_{13} = m/N$; $n_{13} = m$. If we define $T_y^2 = V_y - C_y + S_y; T_x^2 = V_x - C_x + S_x; and$ $\rho = (\nabla_{xy} - C_{xy} + S_{xy}) / (T_x T_y) ,$ then the minimal variance of (24) is: min $[var(\overline{y}'_t)] = C_{yy} - S_{yy}/N + \frac{T^2_y}{n} \cdot \frac{(1+\sqrt{1-\rho^2})}{2}$, (25) a,b,U

where U=u/n

The variance of the simple sample estimate for the second occasion, $\overline{y}_{t} = \frac{1}{n}(m\overline{y}_{mt} + u\overline{y}_{ut})$, is:

$$\operatorname{var}(\overline{y}_{t}) = C_{yy} - S_{yy} / N + T_{y}^{2} / n \quad . \tag{26}$$

The factor $\frac{1}{2}(1+\sqrt{1-\rho^2})$ in (25) represents the reduction in the third term of the variance of the simple estimate, (26), obtained by using matched samples and is the same in form as obtained in classical sampling theory without response errors (e.g. in Cochran [3]). However, if the correlated response error C is large

relative to T_y^2 , the matching will not signifi-cantly reduce the total variance.

4. References:

- [1] Bailar, A. and Dalenius, T. "Estimating the Response Variance Components of the U.S. Bureau of Census' Survey Model." Sankhya B, 31 (1969), 341-60.
- [2] Chai, J.J. "Correlated Measurement Errors and the Least Squares Estimator of the Regression Coefficient." Jour. Amer. Statist. Assoc. 66 (1971), 478-83.
- [3] Cochran, W.G., Sampling Techniques, 2nd edition, Wiley, 1963, 342-4.
- [4] Koch, G.G. "An Alternative Approach to the Multivariate Response Error Models, with Applications to Estimators Involving Subclass Means." (1972). Submitted for publication.
- Hansen, M.H., Hurwitz, W.N. and Bershad, M.A. [5] "Measurement Errors in Censuses and Surveys." Bull. Inter. Statist. Inst. 38: 2, (1961), 359-74.