Robert J. Casady, National Center for Health Statistics

1. Introduction

The balanced repeated replication (BRR) technique as first proposed by McCarthy (1966) and later more fully developed by Kish and Frankel (1970) has been widely adopted for estimating variances for statistics from complex sample surveys. Unfortunately, this technique has not permitted the estimation of the components of variation inherent in a multistage sampling plan. Motivated by the analytical study of a linear first order estimator from a multistage design a BRR technique is developed to estimate these components of variation. An emperical study and a discussion of the proposed technique in the case of both linear and non-linear estimators is also included.

2. The Design and First Order Estimator

Suppose we have a two stage sampling design where x_{ijk} represents the measurement of interest on the kth element from the jth first stage unit in the ith stratum where i=1, 2, ..., L, j=1, 2, ..., N_i and k=1, 2, ..., M_{ij}. Further, let f_i and f_{ij} be the first and second stage sampling fraction respectively and let n_i and m_{ij} be the first and second stage sample sizes respectively. X' is the usual inflation estimator of the population aggrigate X.

Assuming simple random sampling without replacement at both levels it can be shown that Var (X') = $\sum_{i=1}^{L} N_{i}^{2} (1-f_{i}) \sigma_{i}^{2} / n_{i}$ + $\sum_{i=1}^{L} \sum_{i=1}^{N_{i}} N_{i} M_{ij}^{2} (1-f_{ij}) \sigma_{ij}^{2} / n_{i} m_{ij}$ (1)

where
$$\sigma_{ij}^2 = \sum_{k=1}^{M} (X_{ijk} - \overline{X}_{ij.})^2 / (M_{ij} - 1)$$

 $\sigma_i^2 = \sum_{j=1}^{N} M_{ij}^2 (\overline{X}_{ij.} - \overline{X}_{i..})^2 / N_i - 1)$
 $\overline{X}_{ij.} = \sum_{k=1}^{M} x_{ijk} / M_{ij}$
and $\overline{X}_{i..} = (\sum_{j=1}^{N} \overline{X}_{ij.}) / N_i$.

If the sampling is done with replacement then the terms f_i and f_{ij} are set equal to zero in equation (1).

3. The BRR Second Order Estimators

3.1 The Usual BRR Estimator of Variance

Let us assume the stratified two stage design described in the previous section and in addition assume that n_i and m_i are even

to avoid certain tedious but unconsequential algebraic complications. One procedure suggested by Kish and Frankel (1970) for BRR variance estimation was to partition the sample population into two groups by randomly assigning each of the first stage units in each stratum to one of two equal sized groups. Pseudo-replicates are formed by selecting one of the two groups from each of the strata and then each of these pseudoreplicates is used to calculate a numerical value for the first order estimator. Using the above procedure it can be shown that the expectation of

the BRR estimator S_{hs}^2 of the variance of the simple inflation estimator X' is given by

 $E (S_{hs}^{2}) = \sum_{i=1}^{L} N_{i}^{2} \sigma_{i}^{2}/n_{i}$ + $\sum_{i=1}^{L} \sum_{j=1}^{N_{i}} N_{j}M_{ij}^{2} (1-f_{ij}) \sigma_{ij}^{2}/n_{i}m_{ij}$ (2)

when the sampling is without replacement. When the sampling is with replacement the f_{ij} are set equal to zero in equation (2).

It should be noted that when the half-sample estimates are calculated as described above the variation among them is a reflection of both the first and second stage sources of variation in Var (X'), however, they are not in the correct proportion when the sampling is with replacement. Even if one makes simplifying assumptions as in Frankel (1971) it is not possible at this stage to modify the BRR estimator to give a unbiased estimate of Var (X'). Thus if one is interested in such areas as optimum allocation, correcting the usual BRR estimator for bias, etc. it is necessary that an estimator of each of the components of variation be available. In what follows the basic ideas of the BRR method are modified to provide estimates of the components of variation.

3.2 BRR Estimates of Variance Components

T.

First, each of the sampled first stage units is to be considered a pseudo-stratum thus result-

ing in
$$\sum_{i=1}^{n}$$
 n pseudo-strata. Secondly, within i=1

each of the pseudo-strata the second stage units are randomly placed in one of two equal sized groups. By selecting one of the two groups from each of the pseudo-strata the sample elements are partitioned for use in a half-sample estimate. Observe that differences between half-sample estimates reflect only second stage sampling. Denoting the BRR variance estimator formed by using a balanced set of these half samples as

$$S_{hs,2}^{2} \text{ it can be shown that}$$

$$E (S_{hs,2}^{2}) = \sum_{\substack{i=1 \\ i=1 }}^{L} \sum_{\substack{i=1 \\ i=1 }}^{N} N_{ij} \sigma_{ij}^{2} / n_{i} m_{ij}$$

Thus, in the case of sampling with replacement this BRR method can be used to provide an estimate of the component of variation due to second stage sampling and in the case of linear first order estimator the estimator will be unbiased. Further, the estimator

(3)

 $S_{hs,1}^2 = S_{hs}^2 - S_{hs,2}^2$ provides an estimator for

the first-stage component of variance which is also unbiased for linear first order estimators.

When sampling without replacement it is necessary to make simplifying assumptions to get unbiased BRR variance estimators even when dealing with linear first order estimators. Extending Frankel's (1971) assumptions for a one stage design to the two stage design we assume that $f_i = f_1$ and $f_{ij} = f_2$ for all i and j.

It then follows from equations (1), (2) and (3) that for linear estimators

$$(1-f_2)S_{hs,2}^2$$
, $S_{hs,1}^2 = (1-f_1)(S_{hs}^2-(1-f_2)S_{hs,2}^2)$

and $S_{hs,1}^2$ + (1-f₂) $S_{hs,2}^2$ are unbiased estimators

of the second stage component, the first stage component and total variance respectively.

4. A Numerical Example

The following example shows how the method outlined in the previous section can be used to calculate the variance components and the total variance in the case of a none linear estimator.

Suppose we have data from a survey where there are 3 strata with 2 first stage sample units per stratum and 2 second stage units per first stage unit. The sample units are denoted by I_{iik} which represents the kth second stage

unit within the jth first stage unit within the

ith stratum. Assume that

f_i = .5 for i-1, 2, 3

and $f_{ij} = .1$ for i=1, 2, 3 and j=1, 2.

Now suppose that two characteristics are measured on each sample unit, namely characteristic X and characteristic Y. We will let x and y represent the measurements of ijk

X and Y respectively on I _____ijk. The combined ratio estimator

will be used to estimate R= E(X)/E(Y). The (hypothetical) data is summarized below

$x_{111} = 10$	$y_{111} = 30$
$x_{112} = 12$	$y_{112} = 38$
$x_{121} = 21$	$y_{121} = 50$
$x_{122} = 17$	$y_{122} = 52$
$x_{211} = 8$	$y_{211} = 16$
$x_{212} = 11$	$y_{212} = 20$
$x_{221} = 9$	$y_{221} = 30$
$x_{222} = 12$	$y_{222} = 48$
$x_{311} = 14$	y ₃₁₁ = 50
$x_{312} = 20$	$y_{312} = 72$
$x_{321} = 20$	$y_{321} = 42$
$x_{322} = 24$	$y_{322} = 49$

It can be verified that $\hat{R} = 2.7921$ and to estimate Var (\hat{R}) we must calculate S^2_{hs} and $S^2_{hs,2}$.

First to calculate S_{hs}^2 the sample population is grouped as indicated below

Strata	Group I	Group II		
1	(I ₁₁₁ I ₁₁₂)	(I ₁₂₁ I ₁₂₂)		
2	(I ₂₁₁ I ₂₁₂)	(I ₂₂₁ I ₂₂₂)		
3	(I ₃₁₁ I ₃₁₂)	(I ₃₂₁ I ₃₂₂)		

The selection matrix

is then used to form the four half samples which produce the $\hat{R}_1 = 3.0133$, $\hat{R}_2 = 2.2673$, $\hat{R}_3 = 2.7241$ and $R_4 = 3.2473$. Using these half sample estimates we find that

$$S_{hs}^{2} = \sum_{i=1}^{4} (\hat{R}_{i} - \hat{R})^{2}/4 = .1340.$$

Now to calculate $S^2_{hs,2}$ the data is grouped as indicated below

Pseudo-Strata	Group I	Group II	
1	I ₁₁₁	1 ₁₁₂	
2	1 ₁₂₁	1 ₁₂₂	
3	^I 211	1 ₂₁₂	
4	^I 221	1 ₂₂₂	
5	I ₃₁₁	1 ₃₁₂	
6	^I 321	^I 321	

and then using the selection matrix

1	1	-1	-1	1	1	-1	-1)
1	-1	-1	1	1	-1	-1	1
1	1	1	1	-1	-1	-1	-1
1	-1	1	-1	-1	1	-1	1
1	1	-1	-1	-1	-1	1	1
\backslash_1	-1	-1	1	-1	1	1	-1)

we form eight half samples which produce the $\ensuremath{\mathsf{estimates}}$

$$R_1 = 2.6585, R_2 = 2.8824, R_3 = 2.8556,$$

 $\hat{R}_4 = 2.8602, \hat{R}_5 = 2.7449, \hat{R}_6 = 2.8276,$
 $\hat{R}_7 = 2.9070, \hat{R}_8 = 2.6044$

Using these half-sample estimates we find that

$$S_{hs,2}^2 = \sum_{i=1}^{8} (\hat{R}_i - \hat{R})^2 / 8 = .0108$$

Thus the estimate for the component of variance due to second stage sampling would be

$$(1-f_2) S_{hs,2}^2 = .9(.0108) = .0097).$$

The estimate for the first stage component of variance is

$$(1-f_1)$$
 $(S_{hs}^2 - (1-f_2) S_{hs,2}^2) = .0622$

Then adding the estimates for the two components together we have

 $\hat{Var}(\hat{R}) = .0622 + .0097 - .0719.$ It is of interest to note that prior to this paper S_{hs}^2 , which is 86% larger than $\hat{Var}(\hat{R})$,

would have been used to estimate Var (R).

5. The Monte Carlo Study

5.1 The Populations, Sample Design and First Order Statistics

For the Monte Carlo study of the BRR variance component estimators three populations were generated on the computer. Each of the three populations consisted of 10,000 individuals who were grouped into 10 equal sized strata. Each stratum was in turn partitioned into 10 groups of 100 individuals each. These groups were the first stage sampling units and the individuals within the groups were the second stage sampling units.

Measurements of two characteristics were available for each individual in the population. The numeric values for the two measurements were denoted by x_{ijk} respectively for the

 k^{th} individual in the jth group in the ith stratum. For this study x and y were ijk

realizations of the random variables.

$$X_{ijk} = \mu_{ijk} + (C_{ijk} - 5.0) / \sqrt{2.5}$$

and
$$Y_{ijk} = R X_{ijk} + \varepsilon_{ijk}$$

where $C_{ijk}^{} \sim \chi^2$ (5)

$$\varepsilon_{ijk} \sim N(0, \sigma^2)$$

and C_{ijk} and ε_{ijk} were independent for all i, j and k. The values of the μ_{ij} 's may be found in

Table 1. The parameter R was equal to 2 in all of the populations and for Population I, $\sigma = 2$; for Population II, $\sigma = 100$ and for Population III, $\sigma = 17.3205$. These values were chosen so the correlation between x_{iik} and y_{iik} would be .99,

.1 and .5 respectively.

The sample design was a stratified two stage design using simple random sampling without replacement at both levels. Two first stage units were drawn from each stratum and four second stage units were selected from each first stage unit. From the sample three first order statistics were calculated. They were the usual inflation estimator of the aggregate of the x-measurements

$$x' = \frac{10,000}{80} \sum_{i=1}^{10} \sum_{j=1}^{2} x_{ijk},$$

the usual inflation estimator of the aggregate of the y-measurements

$$y' = \frac{10,000}{80} \sum_{i=1}^{10} \sum_{j=1}^{2} \sum_{k=1}^{4} y_{ijk}$$

and the combined ratio estimator of the aggregates

$$R = y' / x'.$$

5.2 Second Order Estimators

The variance of each of the first order estimators and the covariance of x and y were calculated by the following methods:

<u>Table 1</u> - The values of the parameter $\boldsymbol{\mu}_{\textbf{ij}}$

					STRATA					
1st STAGE UNITS	1	2	3	4	5	6	7	8	9	10
1.	35.5	40.5	45.5	50.5	55.5	60.5	65.5	70.5	75.5	80.5
2.	36.5	41.5	46.5	51.5	56.5	61.5	66.5	71.5	76.5	81.5
3.	37.5	42.5	47.5	52.5	57.5	62.5	67.5	72.5	77.5	82.5
4.	38.5	43.5	48.5	53.5	58.5	63.5	68.5	73.5	78.5	83.5
5.	39.5	44.5	49.5	54.5	59.5	64.5	69.5	74.5	79.5	84.5
6.	40.5	45.5	50.5	55.5	60.5	65.5	70.5	75.5	80.5	85.5
7.	41.5	46.5	51.5	56.5	61.5	66.5	71.5	76.5	81.5	86.5
8.	42.5	47.5	52.5	57.5	62.5	67.5	72.5	77.5	82.5	87.5
9.	43.5	48.5	53.5	58.5	63.5	68.5	73.5	78.5	83.5	88.5
10.	44.5	49.5	54.5	59.5	64.5	69.5	74.5	79.5	84.5	89.5

Table 2 -	A comparison	of the variance	estimators fo	or Population I	when the	design is stratified
	two stage wit	h random sampli	ng without rep	lacement at bo	oth stages	

		Emperical H	Emperical Variances			
Parameter	True Value	Standard Estimator	Uncorrected Half-sample Estimator	lst and 2nd Stage Corrected Half-sample Estimator	Standard Estimator	lst and 2nd Stage Corrected Half-sample Estimator
Variance (X')	.66031x10 ⁸	.64864x10 ⁸	.73961x10 ⁸	.64862x10 ⁸	.56946x10 ¹⁵	$.56925 \times 10^{15}$
Between Component	.37472x10 ⁸	.36392x10 ⁸		.36399x10 ⁸	.56722x10 ¹⁵	.61161x10 ¹⁵
Within Component	.28559x10 ⁸	.28472x10 ⁸	-	.28463x10 ⁸	.04346x10 ¹⁵	.09788x10 ¹⁵
Variance (Y')	.26758x10 ⁹	.26287x10 ⁹	.29901x10 ⁹	.26286x10 ⁹	.92198x10 ¹⁶	.92177x10 ¹⁶
Between Component	.14911x10 ⁹	.14457x10 ⁹	-	.14461x10 ⁹	.92326x10 ¹⁶	.99753x10 ¹⁶
Within Component	.11847x10 ⁹	.11830x10 ⁹	-	.11825x10 ⁹	.07172x10 ¹⁶	.16350x10 ¹⁶
Variance (R)	.12726x10 ⁻⁴	$.12969 \times 10^{-4}$.13134x10 ⁻⁴	.12965x10 ⁻⁴	.21430x10 ⁻¹⁰	.21904x10 ⁻¹⁰
Between Component	.00423x10 ⁻⁴	.00609x10 ⁻⁴	-	.00675x10 ⁻⁴	$.23547 \times 10^{-10}$.31291×10 ⁻¹⁰
Within Component	.12303x10 ⁻⁴	.12361x10 ⁻⁴	-	.12290x10 ⁻⁴	.04631x10 ⁻¹⁰	$.16357 \times 10^{-10}$
Covariance (X',Y')	.13168x10 ⁹	.12932x10 ⁹	.14743x10 ⁹	.12931x10 ⁹	.22700x10 ¹⁶	.22689x10 ¹⁶
Between Component	.07471x10 ⁹	.07247x10 ⁹	_	.07249x10 ⁹	.22652x10 ¹⁶	.24437x10 ¹⁶
Within Component	.05698x10 ⁹	.05684x10 ⁹	_	.05683x10 ⁹	.01743x10 ¹⁶	.03919x10 ¹⁶

•

		Emperical	Variances			
Parameter	True Value	Standard Estimator	Uncorrected Half-sample Estimator	lst and 2nd Stage Corrected Half-sample Estimator	Standard Estimator	lst and 2nd Stage Corrected Half-sample Estimator
Variance (X')	.67707x10 ⁸	.67733x10 ⁸	.77526x10 ⁸	.67748x10 ⁸	.60799x10 ¹⁵	.61002x10 ¹⁵
Between Component	.39094x10 ⁸	.39175x10 ⁸	_	.39113x10 ⁸	.62000x10 ¹⁵	.66765x10 ¹⁵
Within Component	.28613x10 ⁸	.28558x10 ⁸	_	.28635x10 ⁸	.04036x10 ¹⁵	.11004x10 ¹⁵
Variance (Y')	.13176x10 ¹¹	.13323x10 ¹¹	.13441x10 ¹¹	.13329x10 ¹¹	.23914x10 ²⁰	.24312 x 10 ²⁰
Between Component	.00451x10 ¹¹	.00472x10 ¹¹	_	.00448x10 ¹¹	.28146x10 ²⁰	.38104x10 ²⁰
Within Component	.12726x10 ¹¹	.12850x10 ¹¹	_	.12881x10 ¹¹	.05351x10 ²⁰	.19794x10 ²⁰
Variance (R)	.33143x10 ⁻¹	.33456x10 ⁻¹	.33684x10 ⁻¹	.33500x10 ⁻¹	.15205x10 ⁻³	.15466x10 ⁻³
Between Component	.00849x10 ⁻¹	.00801x10 ⁻¹	_	.00736x10 ⁻¹	.17848x10 ⁻³	.24482x10 ⁻³
Within Component	$.32294 \times 10^{-1}$	$.32655 \times 10^{-1}$	-	.32764x10 ⁻¹	$.03530 \times 10^{-3}$.13004x10 ⁻³
Covariance (X',Y')	.12639x10 ⁹	.13492x10 ⁹	.15505x10 ⁹	.13446x10 ⁹	.73779x10 ¹⁷	.74074x10 ¹⁷
Between Component	.06896x10 ⁹	.08051x10 ⁹	-	.08236x10 ⁹	.79886x10 ¹⁷	.90431x10 ¹⁷
Within Component	.05743x10 ⁹	.05441x10 ⁹	_	.05211x10 ⁹	.06102x10 ¹⁷	.20750x10 ¹⁷

 $\underline{ Table \ 3} \ - \ A \ comparison \ of \ the \ variance \ estimators \ for \ Population \ II \ when \ the \ design \ is \ stratified \ two \ stage \ with \ random \ sampling \ without \ replacement \ at \ both \ stages$

<u>Table 4</u> - A comparison of the variance estimators for Population III when the design is stratified two stage with random sampling without replacement at both stages

		Emperical	Emperical	Variances		
Parameter	True Value	Standard Estimator	Uncorrected Half-sample Estimator	lst and 2nd Stage Corrected Half-sample Estimator	Standard Estimator	lst and 2nd Stage Corrected Half-sample Estimator
Variance (X')	.66963x10 ⁸	.66751x10 ⁸	.75928x10 ⁸	.66715x10 ⁸	.66280x10 ¹⁵	.66331x10 ¹⁵
Between Component	.36933x10 ⁸	.36705x10 ⁸	-	.36852x10 ⁸	.65507x10 ¹⁵	.71514x10 ¹⁵
Within Component	.30030x10 ⁸	.30046x10 ⁸		.29862x10 ⁸	.05168x10 ¹⁵	.12930x10 ¹⁵
Variance (Y')	.62834x10 ⁹	.62329x10 ⁹	.66076x10 ⁹	.62278x10 ⁹	.58486x10 ¹⁷	.59209x10 ¹⁷
Between Component	.15591x10 ⁹	.14987x10 ⁹	_	.15190x10 ⁹	.61057x10 ¹⁷	.71967x10 ¹⁷
Within Component	.47243x10 ⁹	.47342x10 ⁹	_	.47088x10 ⁹	.08036x10 ¹⁷	.25288x10 ¹⁷
Variance (R)	.93790x10 ⁻³	.93019x10 ⁻³	.93389x10 ⁻³	.93068x10 ⁻³	.11931x10 ⁻⁶	.12370x10 ⁻⁶
Between Component	.02878x10 ⁻³	$.01379 \times 10^{-3}$	-	.01282x10 ⁻³	.13963x10 ⁻⁶	.17550×10 ⁻⁶
Within Component	.90912x10 ⁻³	.91640x10 ⁻³	_	.91785x10 ⁻³	.02866x10 ⁻⁶	.09529x10 ⁻⁶
Covariance (X',Y')	.13251x10 ⁹	.13183x10 ⁹	.15004x10 ⁹	.13165x10 ⁹	.45349x10 ¹⁶	.45417x10 ¹⁶
Between Component	.07310x10 ⁹	.07282x10 ⁹	_	.07357x10 ⁹	.45642x10 ¹⁶	.51028x10 ¹⁶
Within Component	.05941x10 ⁹	.05900x10 ⁹	-	.05808x10 ⁹	.03784x10 ¹⁶	.10857x10 ¹⁶

- (a) The standard textbook formulas which give unbiased estimates for Var (x[']), Var (Y[']) and Cov (x['], y[']). These are nearly unbiased for Var (R[°]);
- (b) The balanced half-sample method without correcting for either the first or second stage sampling fraction; (See Casady (1975) for a discussion of BRR estimates of covariance)
- (c) The balanced half-sample method corrected for both levels of sampling.

In addition the "within" and "between" components of variance for Var(x'), Var(y') and $Var(\hat{R})$ and the "within" and "between" components of covariance are calculated by both the standard textbook formulas and by the balanced half-sample method.

5.3 Expectations and Variances of Second Order Estimators

The expectations and variances of all the second order estimators for each of the three populations were estimated by drawing 1000 samples from each of the three populations. These results are given Tables 2, 3, and 4. This study was primarily exploratory and a far more extensive study is necessary before any definitive conclusions can be drawn, however, it would appear that the BRR estimates of the variance components will be satisfactory. Also, the extensions of the BRR method to the estimation of covariance seems to be satisfactory. The most encouraging result at this time is the fact that the variance of the BRR variance estimators are of the same magnitude as the standard estimators.

6. References

- Casady, Robert J. (1975) "A BRR Estimator of Covariance". National Center for Health Statistics (Internal Memorandum).
- Frankel, Martin R. (1971) <u>Inference to Survey</u> <u>Samples An Emperical Investigation.</u> University of Michigan, Ann Arbor, Michigan.
- Kish, Leslie and Frankel, Martin R. (1971) "Balanced Repeated Replication for Standard Errors". Journal of The American Statistical Association, 65, 1071-94.
- McCarthy, Philip J. (1966) "Replication: An Approach to the Analysis of Data from Complex Surveys.": <u>Vital and Health Statistics</u>. PHS Publication No. 1000, Series 2, No. 14. Government Printing Office, Washington, D.C.