

COMPARISONS OF VARIANCE ESTIMATORS IN STRATIFIED RANDOM AND SYSTEMATIC SAMPLING

Richard Valliant, U.S. Bureau of Labor Statistics
Room 2126, 441 G Street N.W., Washington D.C. 20212

Keywords : Jackknife variance estimator; separate ratio estimator; separate regression estimator; superpopulation model.

1. INTRODUCTION

Ratio and regression estimation in conjunction with stratification are familiar and well-studied methods in the survey sampling literature. Design-based variance estimators are summarized by Cochran (1977). Wu (1985) introduced a class of estimators, which included the standard ones, for the combined ratio estimator and obtained the member of the class optimal in terms of design mean squared error (*mse*). In the unstratified case, design-based studies of the ratio estimator have been done by Rao and Rao (1971), Wu (1982), and Wu and Deng (1983). Deng and Wu (1987) also studied design-based properties of variance estimators for the unstratified regression estimator. Conditional model-based studies have been done by Royall and Eberhardt (1975) and Royall and Cumberland (1981a, 1981b) and have been extended to stratification by Valliant (1987a).

Most previous studies have been done in the context of simple random sampling (*srs*) or stratified simple random sampling (*stsr*s) with relatively little attention given to stratified systematic sampling (*stsys*) in the ratio estimation problem. Much of the literature on variance estimation in systematic sampling deals only with the simple sample mean (e.g. Heilbron 1978, Wolter 1984). Iachan (1982) gives an extensive review of studies on systematic sampling and notes that there is a need for work on more complex estimators. This paper contrasts the effects of *stsr*s and *stsys* on properties of variance estimators for ratio and regression estimators. Kott (1986) noted that systematic sampling is one method of protecting against certain kinds of model biases when estimating a mean. As illustrated here, systematic sampling can also have important effects on variance estimators.

The population is divided into H , a fixed number, of strata and within stratum h a sample of n_h units is selected from the total of N_h units. The sampling fraction in stratum h is $f_h = n_h/N_h$ and the set of sample units from stratum h is denoted as s_h . The total population size is $N = \sum_h N_h$ and the total sample size is $n = \sum_h n_h$. The proportion of the population in stratum h is $W_h = N_h/N$. Associated with unit (hi) is a random variable y_{hi} and an auxiliary x_{hi} with the latter known and positive for every unit in the population. Assume that there are bounds B_1 and B_2 such that $0 < B_1 \leq x_{hi} \leq B_2 < \infty$ for each h and i . As in Valliant (1987a,b), for model-based analyses we will consider a situation in which $N_h, n_h \rightarrow \infty, f_h \rightarrow 0$, and n_h/n and W_h converge to constants in all strata.

The finite population means of y and x are $\bar{y} = \sum_h \sum_i y_{hi}/N$ and $\bar{x} = \sum_h \sum_i x_{hi}/N$ and the stratum means are $\bar{y}_h = \sum_i y_{hi}/N_h$ and $\bar{x}_h = \sum_i x_{hi}/N_h$. The separate and combined ratio estimators are defined as

$$\bar{y}_{RS} = \sum_h W_h \bar{y}_{hs} \bar{x}_h / \bar{x}_{hs} \text{ and}$$

$$\bar{y}_{RC} = \bar{y}_s \bar{x} / \bar{x}_s$$

where $\bar{y}_{hs} = \sum_{s_h} y_{hi}/n_h, \bar{x}_{hs} = \sum_{s_h} x_{hi}/n_h, \bar{y}_s$ is the stratified expansion estimator defined as $\bar{y}_s = \sum_h W_h \bar{y}_{hs}$,

and $\bar{x}_s = \sum_h W_h \bar{x}_{hs}$. The separate and combined regression estimators are

$$\bar{y}_{LS} = \sum_h W_h [\bar{y}_{hs} + b_{hs}(\bar{x}_h - \bar{x}_{hs})]$$

$$\bar{y}_{LC} = \bar{y}_s + b(\bar{x} - \bar{x}_s)$$

where $b_{hs} = s_{xyhs}/s_{xxhs}$ and $b = \sum_h K_{1h} s_{xyhs} / \sum_h K_{1h} s_{xxhs}$ with $K_{1h} = W_h^2(1-f_h)/(n_h(n_h-1)), s_{xyhs} = \sum_{s_h}(x_{hi} - \bar{x}_{hs})y_{hi}$, and $s_{xxhs} = \sum_{s_h}(x_{hi} - \bar{x}_{hs})^2$.

We will study these estimators under some special cases of the model

$$\begin{aligned} y_{hi} &= \alpha_h + \beta_h x_{hi} + \epsilon_{hi} \\ E_{\xi}(\epsilon_{hi}) &= 0, \\ \text{var}_{\xi}(\epsilon_{hi}) &= v_{hi} \end{aligned} \quad (1)$$

with the ϵ_{hi} 's uncorrelated. This model is often reasonable when strata are formed based on the size of x and a more complicated relationship between y and x may be approximated linearly within strata. Such populations are often encountered in surveys of business establishments or institutions such as hospitals conducted by national governments.

2. PROPERTIES OF THE RATIO AND REGRESSION ESTIMATORS

Theoretical properties of the ratio and regression estimators are sketched in this section. In order to make comparisons we employ both model and design-based calculations. Two results are useful in this regard. First, under appropriate conditions, $\sqrt{n_h}(\bar{x}_{hs} - \bar{x}_h)$ converges in distribution under simple random sampling without replacement as $n_h \rightarrow \infty$ (Scott and Wu 1981), i.e.

$(\bar{x}_{hs} - \bar{x}_h) = O_d(n_h^{-.5})$ where O_d denotes probabilistic order with respect to the sample design. The second result is due to Kott (1986) and states that when a systematic sample is selected from a list ordered by x and x is bounded as in Section 1, then $(\bar{x}_{hs} - \bar{x}_h) = O(n_h^{-1})$ with the order being nonprobabilistic. Assuming that n_h/n converges to a constant in each stratum, we have $\bar{x}_h/\bar{x}_{hs} = 1 + O_d(n^{-.5})$ under *stsr*s but $\bar{x}_h/\bar{x}_{hs} = 1 + O(n^{-1})$ under *stsys*. It follows that under *stsr*s $\bar{y}_{RS} = \bar{y}_s + O_d(n^{-.5})$ while $\bar{y}_{RS} = \bar{y}_s + O(n^{-1})$ under *stsys*. These same relationships to the stratified expansion estimator \bar{y}_s also hold for $\bar{y}_{RC}, \bar{y}_{LS}$, and \bar{y}_{LC} . Thus, the differences among the four estimators are of small consequence in large systematic samples.

Turning to the model bias and variance of \bar{y}_{RS} under (1), Valliant (1987a) noted that

$$E_{\xi}(\bar{y}_{RS} - \bar{y}) = \sum_h W_h \alpha_h (\bar{x}_h - \bar{x}_{hs}) / \bar{x}_{hs} \quad (2)$$

$$\text{var}_{\xi}(\bar{y}_{RS} - \bar{y}) = \sum_h W_h^2 D_{xh}^2 \bar{v}_{hs} / n_h + O(n^{-1}) \quad (3)$$

where $D_{xh} = \bar{x}_h/\bar{x}_{hs}$ and $\bar{v}_{hs} = \sum_{s_h} v_{hi}/n_h$. The model variance has order n^{-1} , assuming \bar{v}_{hs} converges to a

constant as $n_h \rightarrow \infty$. The model bias (2) is a random variable with respect to the sample design. Since, under *stsr*s ($\bar{x}_{hs} - \bar{x}_h$) = $O_d(n^{-5})$, the square of the bias (2) has order n^{-1} under *stsr*s which is the same order as the model variance (3). On the other hand, under *stsys* the square of the bias is order n^{-2} . The results of Kott (1986) on systematic sampling also can be applied more generally when, for example, $E_\xi(y_{hi})$ is a polynomial in x_{hi} .

Thus, when an *stsys* is selected, the dominant term of the model mean squared error is the leading term of (3) with the square of the model bias being asymptotically much less important than under *stsr*s. Similar arguments lead to the same conclusions for the combined ratio and combined regression estimators. The separate regression estimator is model unbiased under (1) as is well known.

The above results on the size of the model bias have important implications for mean squared error estimation. Earlier research on robust model variance estimation, such as Royall and Cumberland (1978), have concentrated on cases in which $E_\xi(y_{hi})$ is correctly specified. Variance estimators were then developed which were robust under the general variance specification given in model (1). The fact that systematic sampling can reduce the importance of the model biases of the ratio estimators and the combined regression estimator under (1) and under more general models means that there may be hope of successfully estimating their model *mse*'s under that sampling plan.

3. VARIANCE ESTIMATORS

The fact that estimating repeated sampling variances from systematic samples may present special problems not encountered with random samples has long been recognized (e.g. Osborne 1942, Cochran 1946, Wolter 1984). These special problems are often not accounted for in practice. Wolter (1985 ch. 7) notes that common practice in applied survey work is to regard a systematic sample as random and estimate design variances using random sampling formulae. In a population with linear trend, computed variances are often considered to be overestimates because the random sampling formulae do not appropriately reflect the effect of the trend which is picked up by systematic selection (see e.g. Hansen, Hurwitz, and Madow 1953, §11.8, Wolter 1984).

A variety of variance estimators have been studied for \bar{y}_{RC} and \bar{y}_{RS} . This paper examines a number of the choices that have been proposed for use under *stsr*s plans with emphasis on contrasting the properties that obtain under stratified simple random and stratified systematic sampling plans. For \bar{y}_{RS} we include

$$v_{RSg} = \sum_h K_{1h} D_{xh}^g \sum_{sh} r_{1hi}^2 \text{ and}$$

$$v_{RSJ} = \sum_h K_{1h} D_{xh}^2 \left\{ \frac{(n_h-1)}{n_h} \right\}^2 \sum_{sh} \left[\frac{r_{1hi}}{1 - k_{1hi}} - \frac{1}{n_h} \sum_{sh} \frac{r_{1hj}}{1 - k_{1hj}} \right]^2$$

where $r_{1hi} = y_{hi} - x_{hi} \bar{y}_{hs} / \bar{x}_{hs}$ and $k_{1hi} = x_{hi} / (n_h \bar{x}_{hs})$. For the combined ratio estimator we consider

$$v_{RCg} = D_x^g \sum_h K_{1h} \sum_{sh} r_{2hi}^2 \text{ and}$$

$$v_{RCJ} = D_x^2 \sum_h \sum_{sh} K_{1h} \left[\frac{r_{2hi}}{1 - k_{2hi}} - \frac{1}{n_h} \sum_{sh} \frac{r_{2hj}}{1 - k_{2hj}} \right]^2$$

where $D_x = \bar{x} / \bar{x}_s$, $r_{2hi} = (y_{hi} - \bar{y}_{hs}) - (\bar{y}_s / \bar{x}_s)(x_{hi} - \bar{x}_{hs})$, $k_{2hi} = N_h(x_{hi} - \bar{x}_{hs}) / \{(n_h-1)N\bar{x}_s\}$.

The estimators v_{RSg} and v_{RCg} define classes studied by Wu (1985) who found values of g that were optimal in the sense of minimizing the approximate design *mse*'s of the variance estimators. For the separate estimators we treat the case of the same value of g in all strata although Wu proposed that g be allowed to vary among strata. Cases of special interest are $g = 0, 1, 2$ which have been studied by a number of authors. The estimators v_{RSJ} and v_{RCJ} are computational forms for the stratified delete-one jackknife estimator whose general form was defined by Jones (1974). For some estimator $\hat{\theta}$ the general form is $v_J = \sum_h (1-f_h) \{ (n_h-1)/n_h \} \sum_{sh} \{ \hat{\theta}_{(hi)} - \hat{\theta}_{(h)} \}^2$ where $\hat{\theta}_{(hi)}$ has the same form as $\hat{\theta}$ but omits the $(hi)^{th}$ sample unit and $\hat{\theta}_{(h)} = \sum \hat{\theta}_{(hi)} / n_h$. Since all x_{hi} are bounded, k_{1hi} and k_{2hi} are both $o(1)$ and it is clear from the computational forms above that v_{RSJ} is asymptotically equivalent to v_{RS2} , and v_{RCJ} is asymptotically equivalent to v_{RC2} . Wu (1985) earlier showed that under *stsr*s v_{RC2} is the closest approximation to v_{RCJ} within the class v_{RCg} . Royall and Cumberland (1978 §6) also showed that the general jackknife v_J is asymptotically equivalent to a variance estimator, denoted as G_1 by them, which was derived to be robust against failure of the variance specification in a linear model.

Variance estimators we consider for the separate regression estimator are in the class

$$v_{LSg} = \sum_h K_{2h} D_{xh}^g \sum_{sh} d_{1hi}^2$$

where $K_{2h} = W_h^2(1 - f_h) / [n_h(n_h - 2)]$ and $d_{1hi} = (y_{hi} - \bar{y}_{hs}) - b_{hs}(x_{hi} - \bar{x}_{hs})$. For the combined regression estimator consider

$$v_{LCg} = D_x^g \sum_h K_{1h} \sum_{sh} d_{2hi}^2$$

where $d_{2hi} = (y_{hi} - \bar{y}_{hs}) - b(x_{hi} - \bar{x}_{hs})$. The classes defined by v_{LSg} and v_{LCg} were studied by Deng and Wu (1987) for the unstratified case and by Wu (1985). In the empirical study we additionally include the jackknife variance estimators for \bar{y}_{LS} and \bar{y}_{LC} .

In the case of the sample mean Wolter (1984) has studied a number of estimators involving contrasts and other functions of the sample y 's which are designed to address the peculiarities produced by systematic samples. The focus here will not be to develop new variance estimators but to study the consequences of the common practice of using random sampling estimators when the sample is actually systematic.

4. PROPERTIES OF VARIANCE ESTIMATORS

First, consider variance estimators for the separate ratio estimator. Since, for a fixed value of g , $D_{xh}^g = 1 + O_d(n^{-5})$ under *stsr*s, we have $v_{RSg} = v_{RS0} + O_d(n^{-1.5})$ under that plan. However, under systematic sampling $D_{xh}^g = 1 + O(n^{-1})$ and $v_{RSg} = v_{RS0} + O_d(n^{-2})$. Thus, the choice of g is of less consequence when an *stsys* plan is used. Under model (1)

$$E_\xi(v_{RSg}) = \sum_h W_h^2 \frac{D_{xh}^g}{n_h} \left[\bar{v}_{hs} + \alpha_h \frac{S_{xxhs}}{n_h \bar{x}_{hs}^2} \right] \quad (4)$$

where \approx denotes "asymptotically equivalent". Recalling (3), v_{RS2} is approximately model unbiased when $\alpha_h = 0$ while other choices of g lead to a bias. When $\alpha_h \neq 0$, all v_{RSg} are biased estimators of the model mse . The bias may be substantial and positive under $stsys$ because systematic sampling from a list sorted by x prevents small values of s_{xxhs} but reduces the importance of the bias (2). This observation is similar to the findings of Royall and Cumberland (1978, §5.2) on the overestimation by certain variance estimators for the unstratified ($H=1$) ratio estimator in balanced samples ($\bar{x}_{hs} = \bar{x}_h$). On the other hand, if y is extremely variable for a given x so that $\bar{v}_{hs} \gg \alpha_h^2 s_{xxhs} / (n_h \bar{x}_{hs}^2)$, then the model bias of v_{RSg} can be negligible under $stsys$.

Similar theory can be worked out for v_{RCg} . An approximation to $E_{\xi}(v_{RCg})$ is given by the righthand side of (4) with \bar{x}_{hs} replaced by \bar{x}_g . Consequently, the same remarks given above on the model bias of v_{RSg} under $stsys$ also apply to v_{RCg} .

Next, consider the regression estimators. Using the approximation $D_{xh}^g \approx 1 - g(\bar{x}_{hs} - \bar{x}_h)/\bar{x}_h$ and results from Valliant (1987a, §3.3), the approximate model bias of v_{LSg} is

$$\text{bias}_{\xi}(v_{LSg}) \approx \sum_h \frac{W_h^2}{n} (\bar{x}_{hs} - \bar{x}_h) \left[-g \frac{\bar{v}_{hs}}{\bar{x}_{hs}} + \frac{\sum_{sh} (x_{hi} - \bar{x}_{hs}) v_{hi}}{s_{xxhs}} \right]$$

which has order $n^{-1.5}$ under $stsys$ but only n^{-2} under $stsys$. Similar findings apply to v_{LCg} if $\beta_h = \beta$ in all strata. However, if the slope parameter is not the same in all strata, v_{LCg} has a model bias of order n^{-1} as do v_{RSg} and v_{RCg} .

5. SIMULATION RESULTS

The earlier theory was tested in a simulation study using six artificial populations. Use of generated rather than real populations has some advantages in allowing certain population parameters to be systematically varied in order to study their effect on estimator performance. In particular, we controlled (1) curvature of the regression of y on x and (2) the conditional variance of y given x . In each of the six populations 2000 (x,y) pairs were generated. Each x was generated as $x = 150 + 600w$ where w was a standardized chi square random variable with six degrees of freedom (df), i.e. $w = (\chi_6^2 - 6)/\sqrt{12}$. Given x , y was generated as

$$y = a + bx + cx^2 + dx^2z$$

where a , b , c , and d were constants and z was a standardized chi square random variable with six df . Values of x were constrained to be in the interval $[1, 1500]$ while y was restricted to $[50, 2500]$. Table 1 lists the parameter values used for each population and Figure 1 shows scatterplots of samples of 200 units from each population. Populations 1 and 2 both have the same specification for $E_{\xi}(y)$ with population 1 having the variance of y proportional to $x^{1.5}$ while population 2 has

$\text{var}_{\xi}(y) \propto x^2$. The remaining populations are similarly paired.

Each population was divided into five strata with $N_h = 400$ ($h=1, \dots, 5$). From each population four sets of 1000 samples were selected: (1) 1000 stratified simple random samples of size $n=25$ ($n_h=5$ for all h), (2) 1000 $stsys$'s of $n=100$ ($n_h=20$), (3) 1000 $stsys$'s of $n=25$ ($n_h=5$), and (4) 1000 $stsys$'s of $n=100$ ($n_h=20$). All simple random samples were selected without replacement and all systematic samples were selected with random starts after sorting units in ascending order on x within each stratum.

Table 1. Parameters used in generating study populations.

Pop'n	b	c	g
1	1.5	0	.75
2	1.5	0	1.00
3	1.8	-0.0008	.75
4	1.8	-0.0008	1.00
5	-0.3	0.0009	.75
6	-0.3	0.0009	1.00

Note: In all six populations $a=100$ and $d=.5$.

Tables 2 and 3 give root mean square errors ($rmse$'s) for the separate ratio and regression estimators and square roots of the averages of their variance estimators over the sets of 1000 samples. Results for the combined estimators are omitted to conserve space. We emphasize unconditional comparisons, i.e. ones over all 1000 samples, because conditional properties under $stsys$ have been examined elsewhere (Valliant 1987a) and because systematic sampling virtually eliminates conditional differences in the estimators studied here.

First, we examine the precision of the estimators of the mean. In the lower variance populations (populations 1,3,5) the separate ratio estimator has a considerably lower $rmse$ at either sample size under systematic sampling than under random sampling, while in the higher variance populations (2,4,6) differences in the $rmse$'s are small under the two sampling plans. When $n=25$, the separate regression estimator is generally more precise for all populations under $stsys$ than under $stsys$. When $n=100$, the $rmse$'s of \bar{y}_{LS} are similar under random and systematic sampling with the exception of population 4 where $stsys$ is actually more precise. Comparing Tables 2 and 3, there are noticeable differences between the $rmse$'s of \bar{y}_{RS} and \bar{y}_{LS} under random sampling, particularly for $n=25$ in the higher variance populations where \bar{y}_{RS} is more precise. However, in the systematic samples the $rmse$'s of the separate ratio and regression estimates are little different, especially at the larger sample size. This is in accord with the theoretical observation in §2 that \bar{y}_{RS} and \bar{y}_{LS} differ from each other only by a term of order n^{-1} under $stsys$.

Square roots of average variance estimates are also presented in Tables 2 and 3. In random sampling each of the choices of v_{RSg} ($g=0,1,2$) are generally moderate to small underestimates at either sample size. The jackknife v_{RSJ} is somewhat of an overestimate in $stsys$. For the regression estimator \bar{y}_{LS} , all v_{LSg} ($g=0,1,2$) are severe

underestimates in *stsr*s at $n=25$ with the problem being less severe but still present at $n=100$. At $n=25$ with *stsr*s the jackknife v_{LSJ} has especially wild behavior, overestimating in all populations with some of the worst cases being the high variance populations 2, 4, and 6. When $n=100$ the jackknife for \bar{y}_{LS} is the best performer in *stsr*s being a slight overestimate in all populations while the other choices tend to be underestimates.

With systematic sampling the picture changes. Differences in performance of the variance estimators are considerably reduced. In Table 2 v_{RSg} ($g=0,1,2$) have virtually the same means in each population as do v_{LSg} ($g=0,1,2$) in Table 3. In the low variance populations 1, 3, and 5 all v_{RSg} are overestimates in *stsr*s at both sample sizes as predicted earlier on the basis of expression (4). On the other hand, in the high variance populations 2, 4, and 6 the pattern of consistent overestimation does not hold. The performance of the v_{LSg} 's is substantially better under *stsr*s than *stsr*s. Their degree of underestimation is reduced or eliminated at $n=25$ and at $n=100$ is relatively minor where present. When $n=100$, the best performer under *stsr*s in terms of bias is v_{LSJ} .

Table 4 gives empirical standard deviations (*s.d.*'s) of the variance estimates. In either random or systematic sampling there are differences in precision among the v_{RSg} and among the v_{LSg} but the differences are of no great consequence. The most dramatic numbers in Table 4 are for the jackknife for separate regression estimator which has enormous *s.d.*'s under *stsr*s with $n=25$, a finding similar to that of Andersson, Forsman, and Wretman (1987) in the context of price index estimation. The potential for high variability of the jackknife was also noted by Wu (1986) in linear model analysis. The extreme variability of the jackknife is reduced by using systematic sampling, particularly for $n=100$.

6. CONCLUSION

In populations where there is a reasonably smooth relationship between a target variable y and an auxiliary x , systematic sampling is a defensive strategy. Systematic sampling within strata protects stratified ratio and regression estimators against certain kinds of model biases by producing samples which are more likely to be balanced on moments of x than are simple random samples. However, that bias protection does not always extend to variance estimators. In some types of populations variance estimators for the separate ratio estimator are subject to severe overestimation in systematic samples which persists even in large samples. In cases in which strata are formed based on the size of x and the regression of y on x can be approximated as a straight line within each stratum, the separate regression estimator is a good choice for controlling model bias. Additionally, in the types of populations studied here, standard variance estimators for the separate ratio estimator perform well in systematic samples as long as stratum sample sizes are moderately large.

NOTE

Any opinions expressed are those of the author and do not reflect policy of the Bureau of Labor Statistics.

REFERENCES

Andersson, C., Forsman, G., and Wretman, J. (1987),

- "Estimating the variance of a complex statistic: a monte carlo study of some approximate techniques," *Journal of Official Statistics*, 3, 251–265.
- Cochran, W.G. (1946), "Relative accuracy of systematic and random samples for a certain class of populations," *Ann. of Math. Statist.*, 17, 164–177.
- (1977), *Sampling Techniques*. New York: John Wiley.
- Deng, L.Y., and Wu, C.F.J. (1987), "Estimation of variance of the regression estimator," *J. Am. Statist. Assoc.*, 82, 568–576.
- Hansen, M.H., Hurwitz, W.N., Madow, W.G. (1953), *Sample Survey Methods and Theory Vol. I*. New York: John Wiley.
- Heilbron, D.C. (1978), "Comparison of estimators of the variance of systematic sampling," *Biometrika*, 65, 429–433.
- Iachan, R. (1982), "Systematic sampling: a critical review," *Int. Statist. Rev.*, 50, 293–303.
- Jones, H.L. (1974), "Jackknife estimation of functions of stratum means," *Biometrika*, 61, 343–348.
- Kott, P.S. (1986), "Some asymptotic results for the systematic and stratified sampling of a finite population," *Biometrika*, 73, 485–491.
- Osborne, J.G. (1942), "Sampling errors of systematic and random surveys of cover-type areas," *J. Am. Statist. Assoc.*, 37, 256–264.
- Rao, P.S.R.S. and Rao, J.N.K. (1971), "Small sample results for ratio estimators," *Biometrika*, 58, 625–630.
- Royall, R.M., and Eberhardt, K.R. (1975), "Variance estimates for the ratio estimator," *Sankhyā C*, 37, 43–52.
- Royall, R.M. and Cumberland, W.G. (1978), "Variance estimation in finite population sampling," *J. Am. Statist. Assoc.*, 73, 351–358.
- (1981a), "An empirical study of the ratio estimator and estimators of its variance," *J. Am. Statist. Assoc.*, 76, 66–77.
- (1981b), "The finite population linear regression estimator and estimators of its variance — an empirical study," *J. Am. Statist. Assoc.*, 76, 924–930.
- Scott, A.J. and Wu, C.F.J. (1981), "On the asymptotic distribution of ratio and regression estimators," *J. Am. Statist. Assoc.*, 76, 98–102.
- Valliant, R. (1987a), "Conditional properties of some estimators in stratified sampling," *J. Am. Statist. Assoc.*, 82, 509–519.
- (1987b), "Some prediction properties of balanced half-sample variance estimators in single-stage sampling," *J. Royal Statist. Soc. B*, 49, 68–81.
- Wolter, K. (1984), "An investigation of some estimators of variance for systematic sampling," *J. Am. Statist. Assoc.*, 79, 781–790.
- (1985), *Introduction to Variance Estimation*. New York: Springer-Verlag.
- Wu, C.F.J. (1982), "Estimation of variance of the ratio estimator," *Biometrika*, 69, 183–189.
- (1985), "Variance estimation for the combined ratio and combined regression estimators," *J. Royal Statist. Soc. B*, 47, 147–154.
- (1986), "Jackknife, bootstrap, and other resampling methods in regression analysis (and rejoinder)," *Ann. of Statist.*, 14, 1261–1295 and 1343–1350.
- Wu, C.F.J., and Deng, L.Y. (1983), "Estimation of variance of the ratio estimator: an empirical study," In *Scientific Inference, Data Analysis, and Robustness*, eds. G.E.P. Box, T. Leonard, and C.F. Wu. New York: Academic Press, 245–277.

Figure 1. Scatterplots of 200 units from each of the 6 simulation study populations.

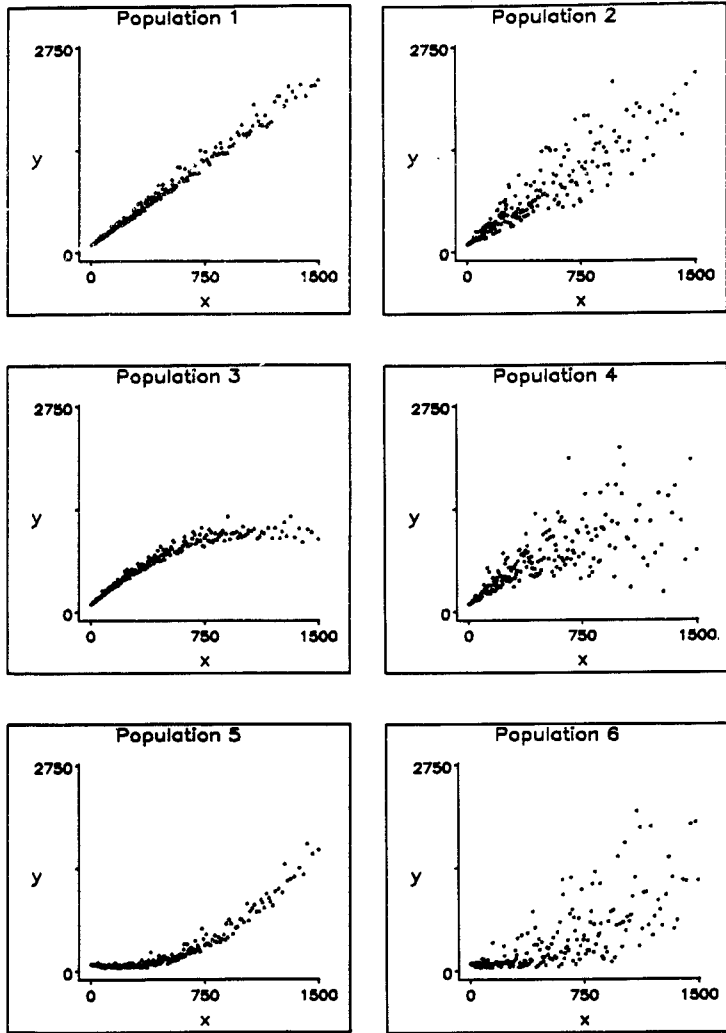


Table 2. Root mean square errors and square roots of average variance estimates for the separate ratio estimator in sets of 1000 stratified simple random and systematic samples from 6 populations.

Pop'n	Sample Type	n	rmse	Square roots of avg. var. ests. in 1000 samples			
				v_{RS0}	v_{RS1}	v_{RS2}	v_{RSJ}
1	ran	25	14.7	13.3	13.6	14.2*	15.3*
		100	6.6	6.4*	6.5*	6.5*	6.5*
	sys	25	11.0	13.7	13.8	13.9	14.0
		100	5.2	6.5	6.5	6.5	6.5
2	ran	25	48.8	48.9*	49.1*	49.4*	50.2*
		100	24.3	24.2*	24.2*	24.3*	24.3*
	sys	25	51.8	48.3	48.3	48.4	48.8
		100	22.9	24.2	24.2	24.2	24.2
3	ran	25	22.5	21.7*	21.8*	22.2*	23.6*
		100	10.6	10.7*	10.7*	10.7*	10.8*
	sys	25	13.7	22.9	22.9	22.9	23.5
		100	5.8	10.9	10.9	10.9	10.9
4	ran	25	60.6	59.2*	59.2*	59.4*	60.3*
		100	29.0	29.3*	29.3*	29.3*	29.4*
	sys	25	59.3	59.7*	59.7*	59.7*	60.3*
		100	34.2	29.1	29.1	29.1	29.1
5	ran	25	20.1	19.9*	20.1*	20.4*	21.8
		100	9.9	9.9*	9.9*	9.9*	9.9*
	sys	25	14.2	21.3	21.3	21.3	21.7
		100	6.7	10.0	10.0	10.0	10.0
6	ran	25	57.9	56.6*	56.6*	56.8*	57.5*
		100	27.2	27.9*	27.9*	27.9*	28.0*
	sys	25	60.9	55.7	55.7	55.8	56.2
		100	28.4	27.9*	27.9*	27.9*	27.9*

*Cases in which the statistic $t =$

$$\left[\bar{v} - \sum_1^S (\bar{y}_{RSi} - \bar{y}) / S \right] / \left\{ \left[\sum_1^S ((v_i - (\bar{y}_{RSi} - \bar{y}))^2) - [\bar{v} - \sum_1^S (\bar{y}_{RSi} - \bar{y})]^2 / S(S-1) \right] \right\}^{-.5}$$

is less than 1.96 in absolute value; $S = 1000$ samples.

Table 3. Root mean square errors and square roots of average variance estimates for the separate linear regression estimator in sets of 1000 stratified simple random and systematic samples from 6 populations.

Pop'n	Sample Type	n	rmse	Square roots of avg. var. ests. in 1000 samples			
				v _{LS0}	v _{LS1}	v _{LS2}	v _{LSJ}
1	ran	25	14.4	11.7	11.7	11.7	24.3
		100	6.0	5.7	5.7	5.7	6.1*
	sys	25	10.8	12.0	12.0	12.0	14.7
		100	5.0	5.7	5.7	5.7	5.9
2	ran	25	61.1	48.1	48.2	48.3	116.2
		100	24.8	23.7*	23.7*	23.7*	25.3*
	sys	25	53.3	48.5	48.5	48.6	61.2
		100	22.7	23.8	23.8	23.8	24.3
3	ran	25	13.7	10.8	10.7	10.7	23.6
		100	5.5	5.3*	5.3*	5.3*	5.7*
	sys	25	10.9	10.4*	10.4*	10.4*	13.1
		100	5.6	5.3	5.3	5.3	5.5*
4	ran	25	74.1	55.9	55.8	55.9	122.3
		100	28.3	27.6*	27.6*	27.6*	29.6*
	sys	25	59.9	55.5	55.6	55.6	69.9
		100	34.4	27.0	27.0	27.0	27.9
5	ran	25	17.2	12.1	12.1	12.1	25.3
		100	6.0	6.0*	6.0*	6.0*	6.4
	sys	25	11.9	12.3*	12.3*	12.3*	15.1
		100	6.1	5.9*	5.9*	5.9*	6.0*
6	ran	25	77.0	56.1	56.0	56.1	136.3
		100	27.7	27.4*	27.4*	27.4*	29.3
	sys	25	62.0	54.7	54.7	54.7	65.9
		100	28.3	27.3*	27.3*	27.3*	28.2*

*Cases in which the statistic t =

$$[\bar{v} - \sum_1^S (\bar{y}_{LSi} - \bar{y})^2 / S] / \left\{ \left[\sum_1^S [(v_i - (\bar{y}_{LSi} - \bar{y}))^2] - [\bar{v} - \sum_1^S (\bar{y}_{LSi} - \bar{y})^2] \right]^2 / [S(S-1)] \right\}^{-5/2}$$

is less than 1.96 in absolute value; S = 1000 samples.

Table 4. Standard deviations of variance estimates for the separate ratio and separate linear regression estimators in sets of 1000 stratified simple random and systematic samples from 6 populations.

Pop'n	Sample Type	n	Standard deviations in 1000 samples							
			v _{RS0}	v _{RS1}	v _{RS2}	v _{RSJ}	v _{LS0}	v _{LS1}	v _{LS2}	v _{LSJ}
1	ran	25	82.7	89.8	122	231	85.2	85.9	88.0	1387
		100	9.2	9.4	10.0	10.1	8.8	8.8	8.8	12.9
	sys	25	85.2	86.0	87.4	88.7	89.8	89.9	90.4	161
		100	9.0	9.1	9.1	9.1	8.9	9.0	9.0	9.7
2	ran	25	1028	1033	1060	1169	1165	1173	1202	41971
		100	114	114	115	115	109	109	111	152
	sys	25	1033	1026	1023	1034	1182	1182	1185	2081
		100	126	127	128	128	119	120	121	122
3	ran	25	224	223	240	345	70.9	67.7	66.1	1675
		100	24.8	24.5	24.9	25.5	7.3	7.1	7.1	10.8
	sys	25	173	169	167	176	53.3	53.3	53.5	115
		100	17.7	17.7	17.8	18.0	7.4	7.4	7.5	8.2
4	ran	25	1725	1708	1725	1789	1691	1675	1691	31915
		100	181	178	178	179	151	149	149	241
	sys	25	1799	1783	1775	1831	1837	1833	1836	3746
		100	174	174	173	174	156	155	155	161
5	ran	25	193	196	223	621	93.6	92.9	94.0	1330
		100	22.7	22.4	22.4	22.8	11.5	11.5	11.5	17.6
	sys	25	186	178	172	180	108	105	101	137
		100	21.2	21.0	20.7	20.9	9.4	9.4	9.4	9.4
6	ran	25	1548	1543	1562	1627	1748	1743	1761	66736
		100	163	162	163	163	163	162	163	256
	sys	25	1503	1506	1514	1546	1595	1603	1615	2592
		100	157	157	157	157	157	156	156	160