

# Jackknife Variance Estimation for Two Samples after Imputation under Two-Phase Sampling

Jong-Min Kim and Jon E. Anderson

Division of Science and Mathematics  
University of Minnesota at Morris  
Morris, MN, 56267, USA

## Abstract

We propose a jackknife variance estimator for the population average from two, two-phase samples after imputation. The jackknife method has long been used to estimate and reduce bias, but has now become a valuable tool for variance estimation. We apply two different sampling methods, (simple random sampling and stratified random sampling) to derive jackknife variance estimators for the two-sample case after imputation under two-phase sampling.

## 1. Introduction

The jackknife method was designed to estimate the bias of an estimator, but now has become a valuable tool for variance estimation since the work of Tukey (1958). In an infinite population context, Tukey (1958) suggested that each jackknife-replicate estimate might be regarded as an independent and identically distributed random variable, which in turn suggests a very simple variance estimator. In the finite population sampling context, each jackknife replicate deletes one unit and modifies the weights of others. Although the one-sample jackknife method was investigated and widely used by many authors, the research on two-sample jackknife estimators is quite limited. Arvesen (1969) was the first to propose such a jackknife estimator as follows: let  $X_1, \dots, X_{n_1}$  and  $X_1, \dots, X_{n_2}$  be two independent samples of sizes  $n_1$  and  $n_2$  from the population. Suppose, in simple random samples of size  $n_1$  and  $n_2$ ,  $r_1$  and  $r_2$  units respond, and  $m_1$  and  $m_2$  do not respond. The observed  $x$  values in population  $k$  are used as donors to create imputed values for the missing data. For example in sample  $k$ , a simple random sample of  $m_k$  values from the  $r_k$  responders could be used as imputed values for the nonresponders. That is, the imputed value  $x_{ik}^*$  for the  $i^{\text{th}}$  observation in sample  $k$  would be one of the observed  $x$

values for one of the responders  $x_{jk}$  in the set of responders  $s_{rk}$ . The imputed estimator of the population mean of  $X$  using a sample  $k$  is given by

$$\bar{x}_{kI} = \frac{1}{n_k} (r_k \bar{x}_{r_k} + m_k \bar{x}_{m_k}), \quad (1)$$

where  $\bar{x}_{m_k}$  is the mean of the imputed values.

## 2. Variance Estimation Under Two-Phase Sampling

Two-phase sampling or double sampling is often employed when it is relatively cheap to take a large preliminary sample in which an auxiliary variable  $x$ , correlated with a characteristic of interest  $y$  alone is measured. The first-phase sample gives a good estimate,  $\bar{x}'$ , of the population mean,  $\bar{X}$ , while the second-phase subsample in which  $y$  is measured is employed to estimate the population mean  $\bar{Y}$  through ratio estimation using  $\bar{x}'$  and  $\bar{y}$ .

### 2.1. Ratio Estimator in Simple Random Sampling

We now extend the two-phase sampling ideas to two samples. Suppose two, first-phase simple random samples  $s'_1$  of size  $n_1$  and  $s'_2$  of size  $n_2$  are taken without replacement from a population of  $N$  elements. Auxiliary attributes  $x_{1i}, x_{2i}$ , are observed for all elements  $i \in s'_1, i \in s'_2$ . Simple random subsamples  $s_1$  of size  $n_1$  and  $s_2$  of size  $n_2$  are taken without replacement from  $s'_1$  and  $s'_2$ . Ratio estimators of  $\bar{Y}$ , are  $\bar{y}_{r_1} = (\bar{y}_1/\bar{x}_1)\bar{x}'_1 = \hat{R}_1\bar{x}'_1$ ,  $\bar{y}_{r_2} = (\bar{y}_2/\bar{x}_2)\bar{x}'_2 = \hat{R}_2\bar{x}'_2$ , where  $\bar{x}'_1$  and  $\bar{x}'_2$  are the means for the first-phase samples  $s'_1$  and  $s'_2$  and  $(\bar{x}_1, \bar{y}_1)$  and  $(\bar{x}_2, \bar{y}_2)$  are the means for the second-phase samples  $s_1$  and  $s_2$ . These estimators are design-consistent for  $\bar{Y}$  as  $n_1 \rightarrow \infty$  and  $n_2 \rightarrow \infty$ .

In single-phase sampling, jackknife variance estimation is often used when  $\bar{X}$ , is known. Under a model-based

framework, Royall and Eberhardt (1975) have shown that the jackknife variance estimator is asymptotically equivalent to a robust variance estimator. As a result, it performs well both unconditionally and conditionally given the ancillary statistic  $\bar{x}$ . The jackknife method for single-phase sampling is not readily applicable to two-phase sampling since  $y_i$  is not observed if  $i \in s'_k - s_k$  in samples  $k = 1, 2$ . We now obtain a jackknife variance estimator for  $\bar{y}_{rk}$  by recalculating  $\bar{y}_{rk}$  with the  $j$ th element removed for each  $j \in s'_k$  then using the variance of these  $n'_k$  jackknife values,  $\bar{y}_{rk}(-j)$ . Clearly, deleting unit  $j$  will affect  $\bar{x}_1$  and  $\bar{x}_2$  only if  $j \in s_k$  and not if  $j \in s'_k - s_k$ , for sample  $k = 1, 2$  while it will affect  $\bar{x}'_k$  for all  $j \in s'_k$ . Thus, we define

$$\bar{y}_{rk}(-j) = [\bar{y}_k(-j)/\bar{x}_k(-j)]\bar{x}'_k(-j),$$

for all  $j \in s'_k$ , where

$$\bar{x}_k(-j) = \begin{cases} \frac{n_k \bar{x}_k - x_j}{n_k - 1}, & \text{if } j \in s_k, \\ \bar{x}_k, & \text{if } j \in s'_k - s_k \end{cases} \quad (2)$$

$$\bar{y}_k(-j) = \begin{cases} \frac{n_k \bar{y}_k - y_j}{n_k - 1}, & \text{if } j \in s_k, \\ \bar{y}_k, & \text{if } j \in s'_k - s_k \end{cases} \quad (3)$$

and  $\bar{x}'_k(-j) = (n'_k \bar{x}'_k - x_j)/(n'_k - 1)$  for all  $j \in s'_k$ . Now apply the usual jackknife method to  $\bar{y}_{rk}(-j)$  to get

$$v_{J_{rk}} = \frac{n'_k - 1}{n'_k} \sum_{j \in s'_k} [\bar{y}_{rk}(-j) - \bar{y}_{rk}]^2. \quad (4)$$

A linearized version of  $v_{J_{rk}}$  for large  $n_k$  is obtained by noting that

$$\bar{y}_{rk}(-j) - \bar{y}_{rk} = \begin{cases} -\hat{R}_k \left( \frac{x_j - \bar{x}_k}{n_k - 1} \right) - \frac{\bar{x}'_k(-j)}{\bar{x}_k(-j)} \left( \frac{x_j - \hat{R}_k x_j}{n_k - 1} \right), & \text{if } j \in s_k, \\ -\hat{R}_k \left( \frac{x_j - \bar{x}_k}{n_k - 1} \right), & \text{if } j \in s'_k - s_k. \end{cases} \quad (5)$$

The jackknife variance estimator is a weighted average of two estimators, given by

$$\begin{aligned} v_{J_r} &= \frac{n'_1 v_{J_{r1}} + n'_2 v_{J_{r2}}}{n'_1 + n'_2} \\ &= \frac{n'_1 - 1}{n'_1 + n'_2} \sum_{j \in s'_1} [\bar{y}_{r1}(-j) - \bar{y}_{r1}]^2 \\ &+ \frac{n'_2 - 1}{n'_1 + n'_2} \sum_{l \in s'_2} [\bar{y}_{r2}(-l) - \bar{y}_{r2}]^2. \end{aligned} \quad (6)$$

Let's consider Rao and Shao (1992) adjusted jackknife variance estimator using data imputation framework. In one

of our samples,  $k = 1, 2$  we define our estimator of  $\bar{Y}$  as  $\bar{y}_{kI} = (1/n'_k) \sum_{i \in s'_k} y_i^*$ . If the observation is part of the second phase sample in sample  $k$ ,  $s_k$ ,  $y_i^* = y_i$ , because it is observed. If the value of  $y$  is not directly observed because it is part of  $s'_k - s_k$  the value is obtained through ratio imputation as  $y_i^* = (\bar{y}_k/\bar{x}_k)x_i$ .

To use convenient one-phase sample variance formulae, Rao and Sitter (1995) proposed the following device to facilitate the computations. They defined

$$\hat{z}_{ki}(-j) = y_{ki}^* + \left\{ \frac{\bar{y}_k(-j)}{\bar{x}_k(-j)} x_{ki} - \frac{\bar{y}_k}{\bar{x}_k} x_{ki} \right\},$$

for sample  $k = 1, 2$ . Under this formulation,  $\hat{z}_{ki}(-j) = y_{ki}^* = (\bar{y}_k/\bar{x}_k)x_{ki}$  for  $j \in s'_k - s_k$  in sample  $k = 1, 2$ , and  $\hat{z}_{ki}(-j) = (\bar{y}_k(-j)/\bar{x}_k(-j))x_{ki}$  for  $j \in s_k$  in sample  $k = 1, 2$ . We also define the adjusted estimator,

$$\bar{y}_{kI}^a(-j) = \frac{1}{n'_k - 1} \sum_{i=1}^{n'_k} \hat{z}_{ki}(-j),$$

and this helps define the jackknife variance estimator for sample  $k$ ,

$$v_{J_{rk}} = \frac{n'_k - 1}{n'_k} \sum_{j \in s'_k} [\bar{y}_{kI}^a(-j) - \bar{y}_{kI}]^2, \quad (7)$$

where  $\bar{y}_{kI} = \bar{y}_{kr}$  under ratio imputation. The jackknife variance estimator based on adjusted imputed estimators  $\bar{y}_{kI}$ ,  $k = 1, 2$  is a weighted average of two estimators, given by

$$\begin{aligned} v_{J_r}^a &= \frac{n'_1 v_{J_{r1}} + n'_2 v_{J_{r2}}}{n'_1 + n'_2} \\ &= \frac{n'_1 - 1}{n'_1 + n'_2} \sum_{j \in s'_1} [\bar{y}_{1I}^a(-j) - \bar{y}_{1I}]^2 \\ &+ \frac{n'_2 - 1}{n'_1 + n'_2} \sum_{l \in s'_2} [\bar{y}_{2I}^a(-l) - \bar{y}_{2I}]^2. \end{aligned} \quad (8)$$

## 2.2. Ratio Estimator in Stratified Random Sampling

Suppose the population of  $N$  units consists of  $L$  strata such that the  $h$ -th stratum consists of  $N_h$  units and  $\sum_{h=1}^L N_h = N$ . Suppose that an auxiliary variable,  $x$ , closely related to an item  $y$  is observed on all sample units,  $s'_{hk}$  in sample  $k = 1, 2$  for stratum  $h$ . Ratio imputation uses  $y_{hki}^* = (\bar{y}_{hk}/\bar{x}_{hk})x_{hki}$  for  $i \in s'_{hk} - s_{hk}$  where  $\bar{y}_{hk}$  and  $\bar{x}_{hk}$  are the means of  $y$  and  $x$  for the respondents  $s_{hk}$  in stratum  $h$ . Ratio imputation can be motivated by the fact that  $y_{hi}^*$

is the best predictor of unobserved  $y_{hki}$  under the following "ratio" superpopulation model  $\xi$ :

$$E_{\xi}(y_{hki}) = \beta_h x_{hki}, \quad V_{\xi}(y_{hki}) = \sigma_h^2 x_{hki},$$

$$\text{cov}_{\xi}(y_{hki}, y_{hkj}) = 0, \quad \text{for } i \neq j \quad (9)$$

provided that the model also holds for the respondents  $s_{hk}$ ; that is, if selection bias is absent. Note that response probabilities can depend on the  $x_{hki}$ 's. Särndal (1992) has named the above equation an "imputation" model. Under ratio imputation,

$$\bar{y}_{kI} = \sum_h^L W_h \bar{y}_{hkI} = \sum_h^L W_h (\bar{y}_{hk} / \bar{x}_{hk}) \bar{x}'_{hk}, \quad (10)$$

where  $\bar{x}'_{hk}$  is the  $x$  mean for the full sample  $s'_{hk}$  from stratum  $h$ . Under the model (9),  $\bar{y}_{kI}$  is design-model unbiased for  $\bar{Y}$ , provided that the model also holds for the respondents. Also, under a uniform response mechanism within each stratum  $h$ , the estimator (10) has the same properties as the standard two-phase sampling separate ratio estimator. This follows by noting that conditionally, given  $n_{hk}$ ,  $s_{hk}$  is a simple random sample of fixed size  $n_{hk}$  drawn from  $s_{hk}$ . Rao (1996) mentioned that the estimator (10) is approximately design unbiased under uniform response in each stratum, proved that  $n_{hk}$  is large for each  $h$ . Note that strata act as imputation classes in the present context. It is readily seen that  $y^*_{hki}(-hki) = [\bar{y}_{hk}(-hki) / \bar{x}_{hk}(-hki)] x'_{hki}$ , under ratio imputation when  $hki$ th respondent is deleted, where

$$\bar{y}_{hk}(-hki) = [n_{hk} \bar{y}_{hk} - y_{hki}] / (n_{hk} - 1) \quad (11)$$

and

$$\bar{x}_{hk}(-hki) = [n_{hk} \bar{x}_{hk} - x_{hki}] / (n_{hk} - 1), \quad (12)$$

for  $k = 1, 2$ .

To use convenient one-phase sample variance formulae, Rao and Sitter (1995) proposed the following device to facilitate the computations. We define

$$\hat{z}_{hki}(-hki) = y^*_{hki} + \left\{ \frac{\bar{y}_{hk}(-hki)}{\bar{x}_{hk}(-hki)} x_{hki} - \frac{\bar{y}_{hk}}{\bar{x}_{hk}} x_{hki} \right\},$$

for sample  $k = 1, 2$ , stratum  $h$ . Under this formulation,  $\hat{z}_{hki}(-hki) = y^*_{hki} = (\bar{y}_{hk} / \bar{x}_{hk}) x_{hki}$  for  $hki \in s'_{hk} - s_{hk}$  in sample  $k = 1, 2$ , and  $\hat{z}_{hki}(-hki) = (\bar{y}_{hk}(-hki) / \bar{x}_{hk}(-hki)) x_{hki}$  for  $hki \in s_{hk}$  in sample  $k = 1, 2$ , stratum  $h$ . We also define the adjusted estimator,

$$\bar{y}_{hkI}^a(-hki) = \frac{1}{n'_{hk} - 1} \sum_{i=1}^{n'_{hk}} \hat{z}_{hki}(-hki),$$

and this helps define the jackknife variance estimator for sample  $k$ .

Using these values, the jackknife variance estimator is given by

$$v_{Jr}(\bar{y}_{kI}) = \sum_{h=1}^L \frac{n'_{hk} - 1}{n_{hk}} \sum_{j=1}^{n'_{hk}} [\bar{y}_{kI}^a(-hki) - \bar{y}_{kI}]^2. \quad (13)$$

For  $k = 1, 2$ , we now obtain a linearized version of the jackknife variance estimator. This variance estimator is useful with computer programs that use the linearization method of variance estimation. Noting that  $\bar{y}_{kI}^a(-hki) - \bar{y}_{kI} = W_h [\bar{y}_{hkI}^a(-hki) - \bar{y}_{hkI}]$ , where  $\bar{y}_{hkI}^a(-hki)$  is the adjusted imputed estimator of the  $h$ th stratum mean  $\bar{Y}_h$  when  $hki$ th sample unit is deleted, we get

$$v_{Jr}(\bar{y}_{kI}) = \sum_{h=1}^L W_h^2 v_{Jr}(\bar{y}_{hkI})$$

$$= \sum_{h=1}^L W_h^2 \left[ \frac{n'_{hk} - 1}{n_{hk}} \sum_{j=1}^{n'_{hk}} (\bar{y}_{hkI}^a(-hki) - \bar{y}_{hkI})^2 \right].$$

The jackknife variance estimator is a weighted average of two estimators, given by

$$v_{Jrs} = \frac{n'_1 v_{Jr}(\bar{y}_{1I}) + n'_2 v_{Jr}(\bar{y}_{2I})}{n'_1 + n'_2}$$

$$= \frac{n'_1}{n'_1 + n'_2} \left( \sum_{h=1}^L W_h^2 \left[ \frac{n'_{h1} - 1}{n'_{h1}} \sum_{j=1}^{n'_{h1}} (\bar{y}_{h1I}^a(-h1j) - \bar{y}_{h1I})^2 \right] \right)$$

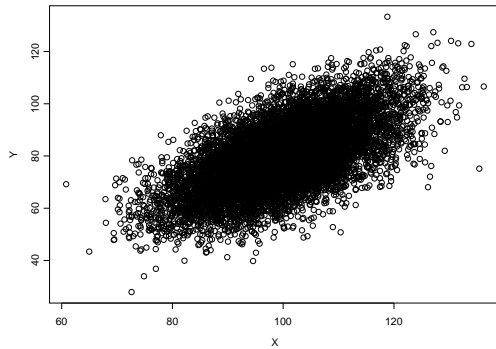
$$+ \frac{n'_2}{n'_1 + n'_2} \left( \sum_{h=1}^L W_h^2 \left[ \frac{n'_{h2} - 1}{n'_{h2}} \sum_{j=1}^{n'_{h2}} (\bar{y}_{h2I}^a(-h2j) - \bar{y}_{h2I})^2 \right] \right)$$

where  $n'_1 = \sum_{h=1}^2 n'_{h1}$  and  $n'_2 = \sum_{h=1}^2 n'_{h2}$ .

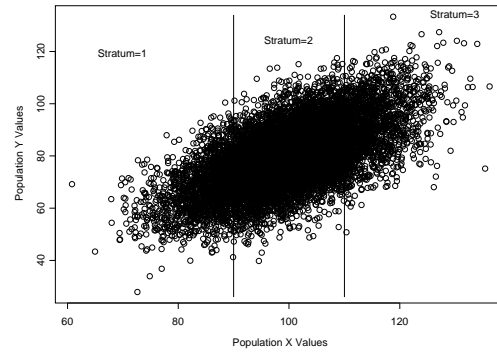
### 3. Simulation Study

For our simulation of the Jackknife variance estimator in simple Random Sampling, we created a population of size  $N = 10000$  by setting  $Y = 0.8 \times X + \epsilon$ , where  $sd(\epsilon) = \sqrt{x}$ . The population is displayed in Figure 3. To study Stratified Random Sampling, we used the same population but we divided the population into three strata:  $X < 90$ ,  $90 \leq X \leq 110$ ,  $110 < X$ . So stratum 1 size is  $N_1 = 1633$  such that  $W_1 = \frac{N_1}{N} = 0.1633$ , stratum 2 size is  $N_2 = 6805$  such that  $W_2 = \frac{N_2}{N} = 0.6805$ , and stratum 3 size is  $N_3 = 1562$  such that  $W_3 = \frac{N_3}{N} = 0.1562$ .

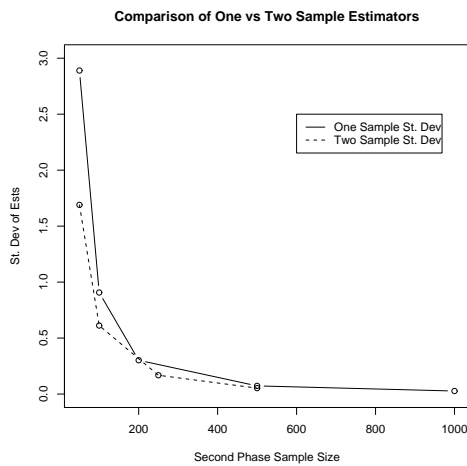
In the figures that follow, we note that there is little difference in the Jackknife variance estimators for the cases studied here, but there are notable differences in the variation of these estimates. Figure 2 shows the reduction in



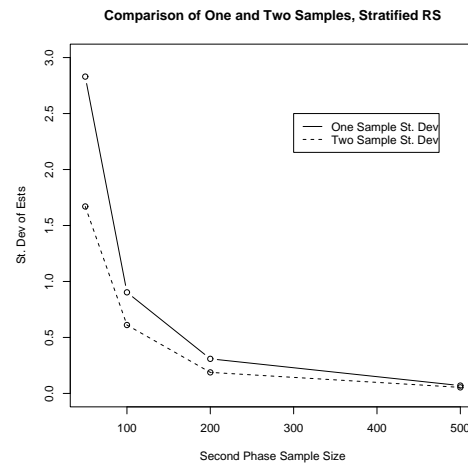
**Figure 1.** *Simulation Study Design: Simple Random Sampling*



**Figure 3.** *Simulation Study Design: Stratified Random Sampling*

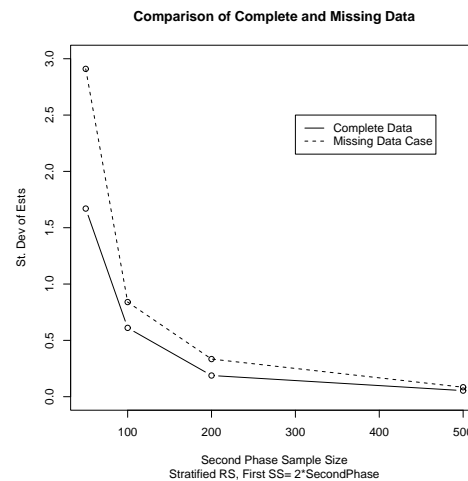


**Figure 2.** *Simple Random Sampling*



**Figure 4.** *Stratified Random Sampling*

variance obtained by increasing the second phase sample size, and how the difference in variance between the one and two sample approaches decreases as the overall study size increases. Figure 4 shows similar behavior for the Jackknife variance estimator for stratified sampling. Both of these figures use a first-phase sample size of 2 times the second phase size. Figure 5 illustrates the increase in variance caused by nonresponse. We note that the increase in variance of our Jackknife variance estimates is only substantial for very small samples, but shows that survey practitioners can gain protection from this instability by using a sufficiently large second phase sample size.



**Figure 5.** *Stratified Random Sampling*

## 4. Conclusions

Figure 2, 4, and 5 show that, in case of small population size, the Jackknife variance estimator for two samples has less standard deviation (variation) than Jackknife variance estimator for one sample. In the future, we intend to study the Jackknife variance estimator for two samples in Stratified Multistage Sampling.

## References

- Arvesen, J.N. (1969). Jackknifing U-statistics. *Ann. Math. Statist.*, **40**, 2076-2100.
- Rao, J. N. K. (1996). On variance estimation with imputed survey data. *Journal of the American Statistical Association*, **91**, 434, 499-506.
- Rao, J. N. K., and Shao, J. (1992). Jackknife variance estimation with survey data under hot deck imputation. *Biometrika*, **79**, 4, 811-822.
- Rao, J. N. K., and Sitter, R. R. (1995). Variance estimation under two-phase sampling with application to imputation for missing data. *Biometrika*, **82**, 2, 453-460.
- Royall, R.M. and Eberhardt, K. R. (1975). Variance estimates for the ratio estimator. *Sankhyā, C*, **37**,1, 43-52.
- Särndal, C.E. (1992). Methods for Estimating the Precision of Survey Estimates when Imputation has been used. *Journal of the American Statistical Association*, **81**, 366-374.
- Tukey, J. (1958). Bias and Confidence in not quite large samples *Ann. Math. Statist.*, **29**, 614.