

## Jackknife variance estimation for functions of Horvitz & Thompson estimators under unequal probability sampling without replacement

Escobar, Emilio L.

*University of Southampton, Social Statistics*

*Southampton, SO17 1BJ, United Kingdom.*

*E-mail: Emilio.Lopez-Escobar@soton.ac.uk*

Berger, Yves G.

*University of Southampton, Southampton Statistical Sciences Research Institute*

*Southampton, SO17 1BJ, United Kingdom.*

*E-mail: Y.G.Berger@soton.ac.uk*

The jackknife is a popular method in survey sampling which is widely used for standard error estimation (e.g. Shao & Tu (1995) and Wolter (2007)). The applicability and theoretical properties of jackknife variance estimators under unequal probabilities without-replacement sampling have been studied to a limited extent. Some examples are given by Campbell (1980), Berger & Skinner (2005), Berger & Rao (2006) and Berger (2007), who proposed jackknife variance estimators for functions of Hájek (1971) point estimators.

We propose two generalised jackknife variance estimators suitable for functions of Horvitz & Thompson (1952) point estimators. Regularity conditions under which the proposed estimators are design-consistent are also provided. These estimators are defined for without-replacement unequal-probability sampling designs and they naturally include finite population corrections. The proposed estimators are compared with the Campbell (1980) jackknife variance estimator for a ratio.

### The class of point estimators

Let  $\mathcal{U} = \{1, \dots, k, l, \dots, N\}$  denote a finite population and let  $s = \{1, \dots, n\}$  denote a sample whose elements are randomly selected with an unequal probability sampling design without replacement,  $s \subseteq \mathcal{U}$ ,  $n \leq N$ . Assume that we are interested in the population parameter  $\theta = h(t_1, \dots, t_q, \dots, t_Q)$  which is a function of population totals from  $Q$  survey variables, where  $h(\cdot)$  is a smooth and differentiable function (e.g. Shao & Tu (1995), Chapter 2),  $t_q = \sum_{k \in \mathcal{U}} y_{qk}$  with  $y_{qk}$  denoting the measurement of the  $q$ -th variable for unit  $k \in \mathcal{U}$ ,  $q = 1, \dots, Q$ . Further, assume we estimate  $\theta$  by the substitution point estimator  $\hat{\theta} = h(\hat{t}_1, \dots, \hat{t}_q, \dots, \hat{t}_Q)$  where  $\hat{t}_q = \sum_{k \in s} w_k y_{qk}$  is the Horvitz & Thompson (1952) point estimator of  $t_q$ , with survey weights  $w_k = \pi_k^{-1}$  where  $\pi_k$  denotes the inclusion probability of unit  $k$ ,  $\pi_k > 0$ ,  $\forall k \in \mathcal{U}$  and  $\pi_{kl}$  denotes the joint inclusion probabilities of units  $k$  and  $l$ ,  $\pi_{kl} > 0$ ,  $\forall k, l \in \mathcal{U}$ .

### The proposed variance estimator

We propose to estimate the variance of  $\hat{\theta}$  by the jackknife variance estimator

$$(1) \quad v_{JHT} = \sum \sum_{(k,l) \in s} \mathcal{D}_{kl} \nu_k \nu_l,$$

with

$$(2) \quad \nu_k = w_k(\hat{\theta} - \hat{\theta}^{(k)}),$$

where  $\mathcal{D}_{kl} = \pi_{kl}^{-1} \{\pi_{kl} - \pi_k \pi_l\}$ , and  $\hat{\theta}^{(k)} = h(\hat{t}_1^{(k)}, \dots, \hat{t}_q^{(k)}, \dots, \hat{t}_Q^{(k)})$  with

$$(3) \quad \hat{t}_q^{(k)} = \sum_{(l \neq k) \in s} w_l y_{ql} + (w_k - 1) y_{qk}.$$

Alternatively, if the sampling design is of fixed sample size, we propose to estimate  $\text{var}(\hat{\theta})$  by

$$(4) \quad v_{JSYG} = -\frac{1}{2} \sum \sum_{(k,l) \in s} \mathcal{D}_{kl} (\nu_k - \nu_l)^2,$$

which is always positive provided  $\mathcal{D}_{kl} < 0$  (e.g. Chao (1982)).

For the simplest case where  $\hat{\theta} = \hat{t} = \sum_{k \in s} w_k y_k$ , (2) and (3) imply  $\nu_k = w_k(\hat{t} - \hat{t}^{(k)}) = w_k(\hat{t} - (\hat{t} - y_k)) = w_k y_k$ . Hence, the proposed jackknives (1) and (4) reduce, respectively, to the Horvitz & Thompson (1952), and the Sen (1953) and Yates & Grundy (1953) unbiased estimators of  $\text{var}(\hat{t})$ .

### Design-consistency

The design-consistency is set under the Isaki and Fuller (1982) asymptotic framework. Accordingly, consider a sequence of nested populations of size  $\{N_t : 0 < N_t < N_{t+1}, \forall t\}$ . Consider also a sequence of (non-necessarily nested) samples of size  $\{n_t : n_t < n_{t+1}; n_t < N_t, \forall t\}$ . Thus,  $t \rightarrow \infty$  implies  $N_t \rightarrow \infty$  and  $n_t \rightarrow \infty$ , with constant  $f = n_t/N_t$ . In what follows, the index  $t$  is dropped to simplify the notation.

In asymptotic studies, it is convenient to work with means instead of totals. Hence, re-define the weights  $w_k$  as  $\tilde{w}_k = w_k/N, \forall k \in \mathcal{U}$ , such that  $\hat{t}_q$  becomes the mean estimator  $\tilde{\mu}_q = \sum_{k \in s} \tilde{w}_k y_{qk}$  for the population mean  $\mu_q = t_q/N, q = 1, \dots, Q$ . Now, recall Results 2.8.1 and 2.8.2 from Särndal *et al.* (1992) and denote by  $\Sigma_{HT} = \sum \sum_{(k,l) \in \mathcal{U}} \mathcal{D}_{kl} \pi_{kl} \tilde{w}_k \tilde{w}_l \mathbf{y}_k \mathbf{y}_l^T, \hat{\Sigma}_{HT} = \sum \sum_{(k,l) \in s} \mathcal{D}_{kl} \tilde{w}_k \tilde{w}_l \mathbf{y}_k \mathbf{y}_l^T, \Sigma_{SYG} = -\frac{1}{2} \sum \sum_{(k,l) \in \mathcal{U}} \mathcal{D}_{kl} \pi_{kl} \{\tilde{w}_k \mathbf{y}_k - \tilde{w}_l \mathbf{y}_l\} \{\tilde{w}_k \mathbf{y}_k - \tilde{w}_l \mathbf{y}_l\}^T$  and  $\hat{\Sigma}_{SYG} = -\frac{1}{2} \sum \sum_{(k,l) \in s} \mathcal{D}_{kl} \{\tilde{w}_k \mathbf{y}_k - \tilde{w}_l \mathbf{y}_l\} \{\tilde{w}_k \mathbf{y}_k - \tilde{w}_l \mathbf{y}_l\}^T$  the multivariate Horvitz-Thompson and Sen-Yates-Grundy population and sample variances of the estimator  $\tilde{\boldsymbol{\mu}} = (\tilde{\mu}_1, \dots, \tilde{\mu}_Q)^T$  for the population parameter  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_Q)^T$ , with  $\pi_{kl} > 0, \forall k, l \in \mathcal{U}$  where  $\mathcal{D}_{kl} = \pi_{kl}^{-1} \{\pi_{kl} - \pi_k \pi_l\}$  and  $\mathbf{y}_k = (y_{1k}, \dots, y_{Qk})^T$ . Consider the following regularity conditions:

- (a)  $v_L/V_L \rightarrow_p 1, V_L \neq 0$  with  $V_L = \nabla(\boldsymbol{\mu})^T \Sigma_{HT} \nabla(\boldsymbol{\mu}), v_L = \nabla(\tilde{\boldsymbol{\mu}})^T \hat{\Sigma}_{HT} \nabla(\tilde{\boldsymbol{\mu}})$  (for fixed sample size designs:  $V_L = \nabla(\boldsymbol{\mu})^T \Sigma_{SYG} \nabla(\boldsymbol{\mu}), v_L = \nabla(\tilde{\boldsymbol{\mu}})^T \hat{\Sigma}_{SYG} \nabla(\tilde{\boldsymbol{\mu}})$ ), where  $\nabla(\mathbf{x})$  is the gradient of  $h(\cdot)$  at  $\mathbf{x} \in \mathbb{R}^Q, \nabla(\mathbf{x}) = (\partial h(\boldsymbol{\mu})/\partial \mu_1, \dots, \partial h(\boldsymbol{\mu})/\partial \mu_Q)_{\boldsymbol{\mu}=\mathbf{x}}^T, h(\cdot)$  is continuous and differentiable at  $\boldsymbol{\mu}$ .
- (b)  $\liminf \{n V_L\} > 0$ .
- (c)  $n^{-1} \sum_{k \in s} \tilde{w}_k^\tau \|\mathbf{y}_k\|^\tau = \mathcal{O}_p(n^{-\tau}), \forall \tau \geq 2$ , with  $\|\mathbf{A}\| = \text{tr}(\mathbf{A}^T \mathbf{A})^{1/2}$  the Euclidean norm.
- (d)  $G_s = n^{-\beta} \sum \sum_{(k \neq l) \in s} (\mathcal{D}_{kl}^-)^2 = \mathcal{O}_p(1)$ , with  $0 \leq \beta < 1, \mathcal{D}_{kl}^- = -\mathcal{D}_{kl}$  if  $\mathcal{D}_{kl} < 0, 0$  otherwise.
- (e)  $H_s = n^{-\beta} \sum \sum_{(k \neq l) \in s} (\mathcal{D}_{kl}^+)^2 = \mathcal{O}_p(1)$ , with  $0 \leq \beta < 1, \mathcal{D}_{kl}^+ = \mathcal{D}_{kl}$  if  $\mathcal{D}_{kl} \geq 0, 0$  otherwise.
- (f)  $\nabla(\mathbf{x})$  is Lipschitz (Hölder) continuous of order  $\delta$ , i.e.  $\|\nabla(\mathbf{x}_1) - \nabla(\mathbf{x}_2)\| \leq \lambda \|\mathbf{x}_1 - \mathbf{x}_2\|^\delta, \lambda > 0$  constant,  $\beta/2 < \delta \leq 1, \mathbf{x}_1$  and  $\mathbf{x}_2$  in neighbourhood of  $\boldsymbol{\mu}$  (e.g. Shao and Tu (1995), page 43).
- (g)  $\|\nabla(\tilde{\boldsymbol{\mu}})\| = \mathcal{O}_p(1)$ .

Condition (a) sets the consistency of the linearisation variance estimator  $v_L$  for  $V_L$  (Särndal *et al.* (1992), Secs. 5.5 & 5.7)). Conditions (b) and (c) are typical (Shao & Tu (1995), pp. 258-260): (b) implies that  $V_L$  decreases with rate  $n^{-1}$  and (c) is a Lyapunov condition. Conditions (d) and (e) are mild requirements on the design, and (f) and (g) are usual smoothness conditions for jackknives.

**Theorem 1.** *For fixed sample size designs, if regularity conditions (a)-(g) hold, then the variance estimator  $v_{JSYG}$  in (4) is asymptotically design-consistent for the approximate linearised variance  $V_L \neq 0$ , i.e.  $v_{JSYG}/V_L \rightarrow_p 1$ .*

**Corollary 1.** *If regularity conditions (a)-(g) hold, then the variance estimator  $v_{JHT}$  in (1) is asymptotically design-consistent for the approximate linearised variance  $V_L \neq 0$ , i.e.  $v_{JHT}/V_L \rightarrow_p 1$ .*

**Corollary 2.** From Theorem 1, by Slutsky's theorem and asymptotic Normality of  $\hat{\theta}$  for  $\theta$ , it follows  $\{v_{JHT}\}^{-1/2}(\hat{\theta} - \theta) \rightarrow_d \mathbf{N}(0, 1)$  and  $\{v_{JSYG}\}^{-1/2}(\hat{\theta} - \theta) \rightarrow_d \mathbf{N}(0, 1)$ . Thus, both jackknife variance estimators,  $v_{JHT}$  and  $v_{JSYG}$  in (1) and (4), allow valid confidence intervals of  $\hat{\theta}$  for  $\theta$ .

**Example: The ratio**

We now illustrate how the proposed estimators works for the ratio point estimator. Let the parameter of interest be  $R = t_y/t_x = \mu_y/\mu_x$ , where  $t_y = \sum_{k \in \mathcal{U}} y_k$  and  $t_x = \sum_{k \in \mathcal{U}} x_k$  are the population totals of the variables  $y$  and  $x$ , and  $\mu_y = t_y/N$  and  $\mu_x = t_x/N$  are population means. Now, assume that  $R$  is estimated with the point estimator  $\hat{R} = \sum_{k \in s} w_k y_k / \sum_{k \in s} w_k x_k$ , which can be thought either as a function of Horvitz-Thompson (1952) total estimators

$$(5) \quad \hat{R} = \hat{t}_y / \hat{t}_x,$$

where  $\hat{t}_y = \sum_{k \in s} w_k y_k$ ,  $\hat{t}_x = \sum_{k \in s} w_k x_k$ , or as a function of Hájek (1971) mean estimators

$$(6) \quad \hat{R} = \check{\mu}_y / \check{\mu}_x,$$

where  $\check{\mu}_y = \hat{t}_y / \hat{N}$  and  $\check{\mu}_x = \hat{t}_x / \hat{N}$ , with  $\hat{N} = \sum_{k \in s} w_k$ . Hence, the proposed variance estimator  $v_{JHT}$  in (1) and the Campbell (1980) jackknife variance estimator, below in (7), are comparable as they estimate the variance of the same point estimator. From Berger & Skinner (2005), Campbell's estimator is defined as

$$(7) \quad v_{JC} = \sum \sum_{(k,l) \in s} \mathcal{D}_{kl} \varepsilon_k \varepsilon_l,$$

where  $\varepsilon_k = (1 - w_k / \hat{N})(\check{\theta} - \check{\theta}^{(k)})$  and  $\check{\theta} = g(\check{\mu}_1, \dots, \check{\mu}_p, \dots, \check{\mu}_P)$  is a function of Hájek (1971) mean estimators from  $P$  variables with  $\check{\mu}_p = \hat{t}_p / \hat{N}$ , and where  $\check{\theta}^{(k)} = g(\check{\mu}_1^{(k)}, \dots, \check{\mu}_p^{(k)}, \dots, \check{\mu}_P^{(k)})$  has the same functional form as  $\check{\theta}$  but using  $\check{\mu}_p^{(k)} = (\hat{t}_p - w_k y_k) / (\hat{N} - w_k)$  instead of  $\check{\mu}_p$ .

It is well known (e.g. Särndal et al. (1992), Result 5.6.2), that the approximate linearised variance of  $\hat{R}$  is given by  $V_L = \sum \sum_{(k,l) \in \mathcal{U}} \mathcal{D}_{kl} \pi_{kl} u_k u_l$  where  $u_k = w_k(y_k - R x_k) / t_x$ . Besides, it is also known that an unbiased estimator of  $V_L$  is given by  $v_L = \sum \sum_{(k,l) \in s} \mathcal{D}_{kl} \check{u}_k \check{u}_l$ , where

$$(8) \quad \check{u}_k = \frac{w_k}{\hat{t}_x} (y_k - \hat{R} x_k).$$

Henceforth, the quantities  $\nu_k$  of the proposed jackknife variance estimator  $v_{JHT}$  in (1) are given by

$$(9) \quad \begin{aligned} \nu_k &= w_k \left( \hat{R} - \frac{\hat{t}_y - y_k}{\hat{t}_x - x_k} \right), \\ &= \frac{w_k}{\hat{t}_x} (y_k - \hat{R} x_k) \left( \frac{\hat{t}_x}{\hat{t}_x - x_k} \right), \\ &= \check{u}_k \left( \frac{\hat{t}_x}{\hat{t}_x - x_k} \right), \end{aligned}$$

whereas the quantities  $\varepsilon_k$  of Campbell's jackknife  $v_{JC}$  in (7) are given by

$$(10) \quad \begin{aligned} \varepsilon_k &= \left( 1 - \frac{w_k}{\hat{N}} \right) \left( \hat{R} - \frac{\hat{t}_y - w_k y_k}{\hat{t}_x - w_k x_k} \right), \\ &= \frac{w_k}{\hat{t}_x} (y_k - \hat{R} x_k) \left( \frac{\hat{t}_x}{\hat{t}_x - w_k x_k} \right) \left( \frac{\hat{N} - w_k}{\hat{N}} \right), \\ &= \check{u}_k \left( \frac{\hat{t}_x}{\hat{t}_x - w_k x_k} \right) \left( \frac{\hat{N} - w_k}{\hat{N}} \right). \end{aligned}$$

It can clearly be seen from (9) that  $\nu_k \doteq \check{u}_k$  if  $(\hat{t}_x - x_k)^{-1}\hat{t}_x \doteq 1$ . On the other hand, from (10) we have that  $\varepsilon_k \doteq \check{u}_k$  if  $\hat{N}^{-1}(\hat{N} - w_k) \doteq 1$  and if  $(\hat{t}_x - w_k x_k)^{-1}\hat{t}_x \doteq 1$ .

Thus, the proposed jackknife estimator is a suitable approximation of the Linearisation variance estimator  $v_L$ . It is more accurate than Campbell's jackknife as (9) is less sensitive to (highly-skewed) weights than (10).

## REFERENCES

- Berger, Y.G. (2007). A jackknife variance estimator for unistage stratified samples with unequal probabilities. *Biometrika*, 94, 953-964.
- Berger, Y.G. & Rao, J.N.K. (2006). Adjusted jackknife for imputation under unequal probability sampling without replacement. *J. R. Statist. Soc. B.*, 68, 531-547.
- Berger, Y.G. & Skinner, C.J. (2005). A jackknife variance estimator for unequal probability sampling. *J. R. Statist. Soc. B.* 67, 1, 79-89.
- Campbell, C. (1980). A different view of finite population estimation. *Proc. Surv. Res. Meth. Sect. Am. Statist. Assoc.* 319-324.
- Chao, M.T. (1982). A general purpose unequal probability sampling plan. *Biometrika*. 69, 3, 653-656.
- Hájek, J. (1971). Comment on a paper by Basu, D. in *Foundations of Statistical Inference* (Godambe, V.P. and Sprott, D.A. eds.). p. 236. Toronto: Holt, Rinehart and Winston.
- Horvitz, D.G. & Thompson, D.J. (1952). A generalization of sampling without replacement from a finite universe. *J. Am. Statist. Assoc.* 47, 663-685.
- Isaki, C.T. & Fuller, W.A. (1982). Survey design under the regression superpopulation model. *J. Am. Statist. Assoc.* 77, 377, 89-96.
- Miller, R.G. (1964). A trustworthy jackknife, *Ann. Math. Statist.* 35, 4, 1594-1605.
- Quenouille, M.H. (1956). Notes on bias in estimation. *Biometrika*. 43, 353-360.
- Robinson, P.M. & Särndal C.E. (1983). Asymptotic properties of the generalized regression estimator in probability sampling. *Sankhyā: The Indian Journal of Statistics, B.* 45, 2, 240-248.
- Särndal, C.-E., Swensson, B. & Wretman, J. (1992). *Model Assisted Survey Sampling*. New York: Springer.
- Sen, A.R. (1953). On the estimate of the variance in sampling with varying probabilities. *J. Indian Soc. Agr. Statist.* 5, 119-127.
- Shao, J. & Tu, D. (1995). *The Jackknife and Bootstrap*. New York: Springer.
- Tukey, J.W. (1958). Bias and confidence in not-quite large samples (abst.). *Ann. Math. Statist.* 29, 2, 614.
- Tillé, Y. (2006). *Sampling Algorithms*. New York: Springer.
- Valliant, R., Dorfman, A.H. & Royall, R.M. (2000). *Finite Population Sampling and Inference: A Prediction Approach*. New York: Wiley.
- Wolter, K.M. (2007). *Introduction to Variance Estimation*. 2nd Ed. New York: Springer.
- Yates, F. & Grundy, P.M. (1953). Selection without replacement from within strata with probability proportional to size. *J. R. Statist. Soc. B.* 15, 253-261.