

An Investigation of Stratified Jackknife Estimators Using Simulated Establishment Data Under an Unequal Probability Sample Design

Philip Steel, Victoria McNerney, John Slanta¹

Abstract

Considering two different sampling schemes (Tille and Pareto), we present the results of a Monte Carlo simulation studying the statistical properties of several variance estimators on a synthetic data set, modeled on establishment data. We test the validity of including a varying finite population correction in the formulation of the stratified jackknife (as done with the Yates-Grundy-Sen estimator), and we compare the effects of direct replication of a rate to a Taylor linearization formulation.

Key Words: jackknife, Tille, Pareto, fpc, unequal probability sampling

1. Introduction

The distribution of manufacturing sector establishments tends to be highly skewed. Consequently, to ensure a representative sample, a stratified, unequal probability sampling scheme (a π -ps or PI-PS design) is preferred to a stratified, simple random sampling scheme, and large units are often included with certainty. To reduce variance and to limit costs, a fixed sample size is preferred as well.

Although unequal probability sampling provides a measure of control over the sample, such strategies present difficulties in variance estimation. Approximate sampling formula variances require all joint inclusion probabilities. This makes computation and storage more complicated, especially when more than two units per strata are selected. These statistics are unbiased only under complete response and for estimates of totals; linearization techniques must be used for non-linear estimators. In contrast, replicate variance estimators require fewer sampling parameters and no linearization formulae.

The purpose of our research is to determine whether a replicate variance estimator can be used for key estimates produced by the U.S. Census Bureau's Quarterly Survey of Plant Capacity (QSPC). The QSPC replaces the annual Plant Capacity Utilization (PCU) Survey and has a very similar design. Like its predecessor, the QSPC has a stratified π -ps design, and the key estimate is the Plant Capacity Utilization Rate.

The combination of survey design and estimator (a smooth statistic) suggests that these data are suited to the stratified jackknife replication variance estimation procedure. We will compare that to the current procedure which takes the approximate sampling formula estimates of the variance and covariance of the rate's totals, and inputs them into a standard Taylor series linearization formula to estimate the variance of the rate. This procedure was used by the PCU and is currently used by the QSPC.

The QSPC sample is currently selected via a Pareto sampling scheme. This replaces the Tillé sampling scheme used for both the PCU and the initial quarters of the QSPC.

¹ This report is released to inform interested parties of ongoing research and to encourage discussion. Any views expressed on statistical or methodological issues are those of the authors and not necessarily those of the U.S. Census Bureau.

Consequently, we consider the statistical properties of all the considered variance estimators under both sampling schemes.

This paper describes the results of a Monte Carlo simulation studying the statistical properties of several variance estimators on a synthetic data set. This data set was modeled from a subset of the industries represented in the PCU historical data. Section 2 describes the Tillé and Pareto sampling methodologies and the associated approximate sampling formulae. The stratified jackknife variance estimators are specified in section 3, and section 4 presents the data synthesis and our Monte Carlo simulation results. We conclude in Section 5 with some observations and recommendations.

2. PCU and QSPC Sample Design

The PCU published industry level estimates of plant capacity utilization rate, defined as the ratio of actual production to full production capability. The PCU survey used a Tillé sampling procedure (Tillé, 1996). In this procedure, units are rejected from the population one by one until the desired sample size is reached. The approximate sampling formula variance of a Horvitz-Thompson estimator under Tillé is:

$$v_{DT}(\hat{Y}) = \sum_{h=1}^L \sum_{i=1}^{n_h} \gamma_{hi} \left[i \sum_{j=1}^i \left(\frac{y_{hj}}{\pi_{hj}} \right)^2 - \left(\sum_{j=1}^i \frac{y_{hj}}{\pi_{hj}} \right)^2 \right]$$

where the data are sorted in ascending order within stratum by π_i (Slanta and Fagan, 1997). There are L sampling strata, indexed by h , and

$$\begin{aligned} \pi_{hi} &= \text{inclusion probability for unit } i \\ \gamma_{hi} &= 0 & i = 1 \\ &= \beta_{hi, h1} - \beta_{h(i+1), h1} & \text{if } 1 < i < n \\ &= \beta_{hi, h1} - 1 & i = n \end{aligned}$$

$$\beta_{hi, h1} = \frac{\pi_{hi} \pi_{h1}}{\pi_{hi, h1}}$$

The variance estimate of the plant capacity utilization rate (\hat{Y}_0) is obtained with the following Taylor linearization variance estimator:

$$(1) \quad v_{DT}^T(\hat{Y}_0) \approx \frac{1}{\hat{T}_2^2} \left[v_{DT}({}_1\hat{T}) + \hat{Y}_0^2 v_{DT}({}_2\hat{T}) - 2\hat{Y}_0 \text{cov}_{DT}({}_1\hat{T}, {}_2\hat{T}) \right]$$

where ${}_1\hat{T}$ is the actual plant utilization, ${}_2\hat{T}$ is the full production capability, $\hat{Y}_0 = \frac{{}_1\hat{T}}{{}_2\hat{T}}$,

and the v_{DT} and cov_{DT} are calculated using the approximate sampling formulae.

The PCU sample was selected towards the end of 2004. The initial frame consisted of manufacturing and publication establishments from the 2002 Economic Census and was stratified by 6 digit NAICS (North American Industry Classification System). To reduce coverage bias, additional strata were added to represent establishments that came into business (were born) after 2002. Births were introduced on an annual basis after the initial cohort. These strata were created from the birth frames identified for the Annual

Survey of Manufactures. Each birth stratum contained representatives from a variety of NAICS codes.

The Tillé sampling procedures used to obtain the PCU sample have been replaced by a Pareto sampling scheme for the QSPC. Pareto sampling is a particular case of order sampling. Like Tillé sampling, Pareto sampling is a π -ps procedure that yields a fixed sample size. However, the Pareto sampling procedure is easier to implement, and the approximate sampling variance formula is quite easy to derive. Both procedures produce similar samples with regards to the distribution of units by size and we hope to confirm that stratified jackknife methods for variance estimation work well on both kinds of samples.

We use the following procedure for selecting a Pareto sample. First, for each unit i within stratum h , compute

$$q_{hi} = \frac{u_{hi}(1 - \pi_{hi})}{(1 - u_{hi})\pi_{hi}} \quad \text{where } u \text{ is } U(0,1).$$

The sample of size n_h consists of n_h units with the smallest q_{hi} within stratum h .

An approximate sampling formula variance of a total obtained by Pareto sampling (Rosén, 1997) is

$$v_{DP}(\hat{Y}) = \sum_{\substack{h \\ C_h > 0}} \frac{n_h}{n_h - 1} \sum_{i=1}^{n_h} (1 - \pi_{hi}) \left(\frac{y_{hi}}{\pi_{hi}} - \frac{B_{hy}}{C_h} \right)^2 = \sum_{\substack{h \\ C_h > 0}} \frac{n_h}{n_h - 1} \left(A_{hy} - \frac{B_{hy}^2}{C_h} \right)$$

$$\text{where } A_{hy} = \sum_{i=1}^{n_h} (1 - \pi_{hi}) \left(\frac{y_{hi}}{\pi_{hi}} \right)^2 \quad B_{hy} = \sum_{i=1}^{n_h} (1 - \pi_{hi}) \frac{y_{hi}}{\pi_{hi}}$$

$$C_{hy} = \sum_{i=1}^{n_h} (1 - \pi_{hi}).$$

The approximate sampling formula covariance is given as

$$\text{cov}_{DP}(\hat{Y}, \hat{X}) = \sum_{\substack{h \\ C_h > 0}} \frac{n_h}{n_h - 1} \sum_{i=1}^{n_h} (1 - \pi_{hi}) \left(\frac{y_{hi}}{\pi_{hi}} - \frac{B_{hy}}{C_h} \right) \left(\frac{x_{hi}}{\pi_{hi}} - \frac{B_{hx}}{C_h} \right) = \sum_{\substack{h \\ C_h > 0}} \frac{n_h}{n_h - 1} \left(A_{hyx} - \frac{B_{hy}B_{hx}}{C_h} \right)$$

$$\text{where } A_{hyx} = \sum_{i=1}^{n_h} (1 - \pi_{hi}) \frac{y_{hi}x_{hi}}{\pi_{hi}^2}.$$

Note that these formulae apply only to Horvitz-Thompson estimates. As with the Tillé samples, we apply Taylor's linearization formula to obtain an approximate variance of the plant capacity utilization rate:

$$(2) \quad v_{DP}^T(\hat{Y}_0) \approx \frac{1}{\hat{T}^2} \left[v_{DP}(\hat{T}_1) + \hat{Y}_0^2 v_{DP}(\hat{T}_2) - 2\hat{Y}_0 \text{cov}_{DP}(\hat{T}_1, \hat{T}_2) \right]$$

where v_{DP} and cov_{DP} are as just defined and \hat{T}_1 , \hat{T}_2 , and \hat{Y}_0 are as they appear in (1). These totals and ratio estimates are calculated for each domain (NAICS).

3. Stratified Jackknife Variance Estimators

The standard stratified jackknife replicate estimate $\hat{T}_{h(i)}$ is created by dropping the i^{th} unit from the sample and reweighting the remaining of the n_h-1 units within stratum h by $n_h/(n_h-1)$. In a stratified simple random sample, the stratified jackknife formula is modified to include the finite population correction factor (if necessary) by incorporating it into the replicate sums-of-squares corresponding to the units' strata. With an unequal probability sample design, this pre-multiplication by a stratum-level constant is not appropriate. Instead, we can multiply each replicate i 's contribution to the sums-of-squares by $(1-\pi_i)$, where π_i is the inclusion probability associated with omitted unit, i . We call this factor a varying finite population correction.

For each sampling scheme and estimator (\hat{Y}) , we estimated the variance with the standard stratified jackknife estimate (3) and two variations:

$$(3) \quad v_J(\hat{Y}) = \sum_h \frac{n_h - 1}{n_h} \sum_{i=1}^{n_h} (1 - \pi_{hi})(\hat{Y}_{h(i)} - \hat{Y}_0)^2$$

where $\hat{Y}_{h(i)}$ is the replicate estimate

\hat{Y}_0 is the full sample estimate,

n_h is the number of units within stratum h , and

$(1 - \pi_{hi})$ is a varying finite population correction for unit i in stratum h .

$$(4) \quad v_{JR}(\hat{Y}) = \sum_h \frac{n_h - 1}{n_h} \sum_{i=1}^{n_h} (1 - \pi_{hi})(\hat{Y}_{h(i)} - \bar{\hat{Y}}_h)^2$$

where $\bar{\hat{Y}}_h$ is the average of the replicates in stratum h ,

$$(5) \quad v_{JU}(\hat{Y}) = \sum_h \frac{n_h - 1}{n_h} \sum_{i=1}^{n_h} (\hat{Y}_{h(i)} - \hat{Y}_0)^2$$

$$\text{When } \hat{Y} \text{ is a ratio, } \hat{Y}_{h(i)} = \frac{\hat{T}_{1h(i)}}{\hat{T}_{2h(i)}}, \hat{Y}_0 = \frac{\hat{T}_1}{\hat{T}_2}, \text{ and } \bar{Y}_h = \frac{\sum_{i=1}^{nh} (1 - \pi_{hi}) \frac{\hat{T}_{1h(i)}}{\hat{T}_{2h(i)}}}{\sum_{i=1}^{nh} (1 - \pi_{hi})}.$$

Note that (3) and (4) differ only in the second term of the sums-of-squares -- it is the full sample estimate in (3) and a weighted average of replicate estimates over the stratum in (4). Equation (4) is algebraically equivalent to the approximate sampling formula variance estimate of a total from a Pareto sample, v_{DP} . Equation (5) is included to assess the necessity of including the varying fpc-correction factor in the variance estimator.

Finally, we compute the Taylor linearization variance estimate using the stratified jackknife variance estimates as input. This estimator mimics our standard procedure for non-linear estimators.

$$(6) \quad v_J^T(\hat{Y}) = \frac{1}{\hat{T}_j^2} \left[v_J(\hat{T}_1) + \hat{Y}_0^2 v_J(\hat{T}_2) - 2\hat{Y}_0 \text{cov}_J(\hat{T}_1, \hat{T}_2) \right]$$

where v_J and cov_J are the variances and covariance computed via the “traditional” stratified jackknife method in (3).

4. Monte Carlo Simulation Study

4.1. Generating An Artificial Population

In order to examine multiple sampling processes and their effect on variances, we produced a simulated population modeled on data selected from the Plant Capacity Utilization survey. For a detailed description of the PCU data and sample design see the publication appendices: <http://www.census.gov/prod/2007pubs/mqc1-06.pdf>. Data for our simulation were obtained from the 2004-2006 survey collections. The PCU surveys hundreds of distinct industries. Despite a detailed classification of industry, some industries display complex distributions of size (capacity). We confined our simulation to those industries for which we had some confidence in our ability to model; that is, those that appeared to be distributed as a pure lognormal.

We used a program developed by the Census Bureau, LogNormSim (McNerney and Adeshiyani, 2006), to generate industry level multivariate lognormal data populations, using noncertainty cases from the selected PCU data as training data. Our multivariate lognormal distributions modeled three variables: full production capability, actual production, and payroll. The first two variables are used in the estimate of the plant capacity utilization rate; the third is a measure of size used to calculate the probability of selection.

For each industry, we evaluated the fit of the simulated population to the training data on seven percentiles, kurtosis, skew, standard deviation, mean, and the correlations between the three variables. We also calculated the average percent difference in all percentiles from the survey to the simulated data and the average percent difference in the correlations. From those without major defect in fit, we selected eleven subpopulations

with attention given to variety and the inclusion of several which modeled industries known to be of particular interest in the final survey publication.

After conducting an empirical analysis of the original PCU data, we decided to construct two additional subpopulations to study particular data phenomena. In industry ‘ZZZZZZ’, we randomly set approximately 18% of the values to 0 in a selected modeled industry. This was the highest reported-zero rate observed in original PCU microdata. In industry ‘YYYYYY’, we added an outlier/influential value. The outlier was a sample case, but received a weight very close to one. In all populations, certainty cases were excluded from the training data used for modeling. Instead, after the 13 simulated subpopulations were created, we added the relevant certainty cases from the PCU survey to our simulated population.

Next, we divided our simulated population into 15 distinct strata: 13 (non-certainty) strata based on frame industry; one birth stratum with cases randomly recruited from those 13 subpopulations based on the proportion of births observed in the PCU survey; and one stratum containing all of the certainty cases. Two of the domains, 327332 and 332722, had a high proportion of births (between 1.5 and 2 percent of cases). Nine domains, including ‘YYYYYY’ and 334518, had no births. It seems likely that our modeling criteria had a tendency to exclude NAICS with active birth processes.

From this population, we drew 10,000 independent samples using the Tillé sampling method and 10,000 independent samples using the Pareto sampling method. In each sample we calculated the variance estimates of the plant capacity utilization ratio using the three stratified jackknife methods (v_J from equation 3, v_{JR} from equation 4, v_{JU} from equation 5), the Taylor approximation using the sampling formula variance (i.e., v_{DT}^T , the former production method that used v_{DT} for Tillé samples (from equation 1) or v_{DP}^T , the current production method that uses v_{DP} for Pareto samples (from equation 2)), and Taylor approximation using v_J (v_J^T , which uses equation 6 with stratified jackknife variances of totals--using equation 3--as inputs).

To assess the performance of each method over repeated samples, we calculated the relative bias and relative stability for each variance estimator (v_s) within each domain using the formulae below (Shao and Tu, 1995, p251):

$$\text{Relative Bias} = \frac{\frac{1}{10000} \sum_{s=1}^{10000} v_s}{MSE(\hat{Y})} - 1$$

$$\text{Relative Stability} = \frac{\left[\frac{1}{10000} \sum_{s=1}^{10000} (v_s - V(\hat{Y}))^2 \right]^{1/2}}{MSE(\hat{Y})}$$

$$\text{where } MSE(\hat{Y}) = \frac{1}{10000} \sum_{s=1}^{10000} (\hat{Y}_s - Y)^2 \text{ and } V(\hat{Y}) = \frac{1}{10000} \sum_{s=1}^{10000} (\hat{Y}_s - \bar{Y})^2,$$

\hat{Y}_s is the sample estimate, \bar{Y} is the mean of all 10,000 sample estimates, and Y is the subpopulation value of the plant capacity utilization rate.

4.2. Results

Tables 1 through 4 present the relative bias and relative stability results for the six variance estimators. The relative bias results presented in tables 1 and 2 can be summarized as follows:

- Under either sampling scheme, the Taylor linearization variance approximations (v_{DT}^T and v_J^T or v_{DP}^T and v_J^T) tend to yield higher absolute relative bias than their fpc-corrected jackknife counterparts (v_J and v_{JR}). In the few populations where the currently implemented Taylor linearization method (v_{DT}^T obtained via (6)) has better properties, the v_{JR} variances provide a close second. Note that the majority of these differences in relative bias are fairly trivial;
- The performance of the v_J estimator appears to be greatly affected by existence of an outlier in the sample (see industry ‘YYYYYY’) but not does not appear to be overly affected by the prevalence of zeros in the data (see industry ‘ZZZZZZ’);
- The v_{JR} replicate estimator appears to be fairly resistant to the presence of an outlier as well as being unaffected by a high prevalence of zeros in the data, as do all three Taylor linearization variance estimators ($v_{DT}^T, v_{DP}^T, v_J^T$);
- The relative bias of the uncorrected replicate variance estimator (v_{JU}) is uniformly poor, with overestimation in all cases. This provides strong evidence of the necessity of including a finite population correction in estimates of variance for both of these fixed size sample designs.

Table 1: Relative Bias Using Tillé Sampling

NAICS	v_J	v_{JR}	v_J^T	v_{DT}^T	v_{JU}
313320	0.12%	-0.72%	-2.38%	-1.21%	61.29%
325131	-1.90%	-3.86%	-11.59%	-9.90%	48.51%
325510	-1.96%	-2.12%	-5.38%	-5.48%	14.93%
327121	1.07%	0.75%	-16.92%	-16.09%	12.74%
327332	6.64%	3.30%	2.65%	-2.12%	47.92%
332721	5.52%	5.21%	4.80%	3.70%	38.84%
332722	5.00%	4.02%	2.85%	2.20%	50.60%
332912	0.05%	-1.31%	-3.40%	-4.84%	52.54%
333414	0.00%	-0.12%	-3.94%	-3.73%	30.19%
334416	0.46%	-0.02%	-2.34%	-2.80%	34.64%
334518	16.24%	-2.74%	11.40%	-4.08%	256.00%
YYYYYY	111.08%	1.84%	104.35%	-6.35%	1880.25%
ZZZZZZ	1.86%	0.52%	0.21%	-0.87%	58.44%

Table 2: Relative Bias Using Pareto Sampling

NAICS	v_J	v_{JR}	v_J^T	v_{DP}^T	v_{JU}
313320	-1.98%	-2.81%	-4.45%	-5.22%	57.98%
325131	-1.42%	-3.40%	-11.12%	-12.74%	49.75%
325510	-0.46%	-0.61%	-3.92%	-4.06%	16.70%
327121	0.88%	0.56%	-17.02%	-17.13%	12.57%
327332	7.42%	4.08%	3.36%	0.18%	48.90%
332721	0.04%	-0.25%	-0.63%	-0.92%	31.65%
332722	0.78%	-0.17%	-1.30%	-2.18%	44.66%
332912	1.87%	0.51%	-1.64%	-2.88%	55.47%
333414	-3.33%	-3.45%	-7.12%	-7.22%	25.90%
334416	-3.83%	-4.29%	-6.50%	-6.91%	29.11%
334518	17.09%	-2.03%	12.22%	-6.13%	261.73%
YYYYYY	106.30%	-0.58%	99.77%	-8.16%	1840.95%
ZZZZZZ	-0.64%	-1.93%	-2.24%	-3.48%	54.57%

The stability results presented in tables 3 and 4 can be summarized as follows:

- In most cases, the stabilities of the variance estimates are comparable for the two corresponding Taylor linearization variance estimators (v_{DT}^T and v_{DP}^T) and the two corrected jackknife variance estimators (v_J and v_{JR}). Thus, direct replication of the plant capacity utilization rate (a ratio) does not appear to detrimentally affect the variance of the variance estimates. It appears that the additional smoothing by the Taylor linearization approximation is not necessary with this ratio estimator, population, and sampling scheme (Tillé or Pareto);
- In the few cases where the stability is affected by the choice of variance estimator, the v_{JR} formulation generally produces the more stable variance estimates than the other methods;
- The uncorrected jackknife replicate estimates (v_{JU}) are the most unstable in all domains.

Table 3: Relative Stability Using Tillé Sampling

NAICS	v_J	v_{JR}	v_J^T	v_{DT}^T	v_{JU}
313320	41.78%	41.54%	39.23%	40.89%	78.13%
325131	69.51%	69.19%	60.44%	64.95%	98.18%
325510	49.95%	49.87%	45.89%	46.01%	56.34%
327121	86.80%	86.39%	66.04%	66.73%	93.43%
327332	41.30%	40.31%	36.31%	35.36%	65.89%
332721	16.30%	16.16%	15.76%	15.39%	42.62%
332722	33.25%	32.86%	31.13%	31.26%	62.57%
332912	30.47%	30.30%	28.68%	29.02%	71.86%
333414	34.48%	34.43%	32.17%	33.03%	49.10%
334416	37.23%	37.05%	35.31%	35.87%	56.52%
334518	70.75%	61.60%	63.88%	62.12%	271.21%
YYYYYY	130.84%	47.66%	123.20%	47.48%	1905.76%
ZZZZZZ	25.84%	25.61%	24.91%	25.12%	66.32%

Table 4: Relative Stability Using Pareto Sampling

NAICS	v_J	v_{JR}	v_J^T	v_{DP}^T	v_{JU}
313320	41.37%	41.18%	38.96%	38.91%	75.03%
325131	70.53%	70.22%	61.09%	61.33%	99.86%
325510	50.83%	50.75%	46.58%	46.53%	57.75%
327121	87.02%	86.61%	66.15%	66.10%	93.62%
327332	42.38%	41.33%	37.15%	36.53%	67.47%
332721	14.44%	14.41%	14.16%	14.14%	35.72%
332722	31.85%	31.57%	30.05%	29.87%	57.07%
332912	30.89%	30.64%	28.86%	28.80%	74.37%
333414	33.29%	33.26%	31.49%	31.48%	45.39%
334416	35.93%	35.83%	34.44%	34.40%	51.79%
334518	71.76%	62.44%	64.81%	57.60%	276.92%
YYYYYY	126.25%	46.67%	118.87%	44.90%	1865.69%
ZZZZZZ	25.22%	25.12%	24.48%	24.48%	62.67%

With either sampling scheme, the v_{JR} estimates have relatively trivial bias and are fairly stable. The estimates v_{JR} and v_J differ significantly only on the outlier population. In all cases, the replicate variance estimates have comparable – if not better – statistical properties to their Taylor linearization counterparts. This presents a clear advantage over the current production method: obtaining estimates with comparable statistical properties in a simpler fashion, with no assumed approximations (e.g., negligible higher order derivatives for the Taylor linearization, trivial non-response for the input approximate sampling formula variances and covariances). These advantages apply regardless of the sampling scheme in this simulation.

5. Conclusion

This research was motivated by a practical problem: to determine whether the stratified jackknife variance estimator could be used to produce usable variance estimates for the U.S. Census Bureau’s Quarterly Survey of Plant Capacity (QSPC). As a part of this analysis, we evaluated both direct replication of a ratio estimate and a variance approximation that combined approximate sampling formula or replicate estimates with a Taylor linearization approximation.

With a highly stratified sample, the literature supports the usage of the stratified jackknife variance estimator for a smooth statistic such as a ratio (e.g., Shao and Tu, 1995). However, we found very little guidance in the literature on the appropriate finite population correction factor with an unequal probability sample. We also found little in terms of concrete recommendations for the “gold standard” in the replicate sums-of-squares. The presented analysis provides strong evidence for the usage of the varying finite-population correction factor in the replicate variance estimator, demonstrating that failing to include this correction leads to consistent overestimation. [Note that Deville (1999) advocates the usage of this adjustment in his paper on linearization variance estimators.] Our research also provides support for the local replicate weighted average in the sums-of-squares, though the performance is only marginally better than the full sample estimate in all but the outlier population.

The strong performance of v_{JR} on the population with an outlier in the current production variable could not be traced to any component of the estimator. It may only be coincident to whatever direction the outlier pulls the ratio. More research is needed to determine if in fact the v_{JR} estimator performs better on data with significant outliers.

Acknowledgements

The authors would like to thank Katherine J. Thompson and David Kinyon for their help in the study design and their comments on this paper.

References

Deville, J.C. (1999). "Variance Estimation for Complex Statistics and Estimators: Linearization and Residual Techniques", *Survey Methodology*, **25**, 2, pp. 193-203.

Tillé, Yves (1996) "An Elimination Procedure for Unequal Probability Sampling Without Replacement", *Biometrika*, **83**, 1, pp. 238-241.

Slanta, John G. and Fagan, James T. (1997), "A Modified Approach to Sample Selection and Variance Estimation with Probability Proportional to Size and Fixed Sample Size", unpublished report, MCD Working Paper Number: Census/MCD/WP-97/02.

McNerney, Victoria G., and Adeshiyan, Samson A. (2006), "User Guide for Generalized Population Simulation Programs", unpublished report.

Rosén, Bengt (1997), "On Sampling with Probability Proportional to Size", *Journal of Statistical Planning and Inference*, **62**, pp.159-191.

Shao, J. and Tu, D. (1995). The Jackknife and Bootstrap. New York: Springer Verlag.