



An Alternative Stratified Three-Stage Estimators for Finite Population Total

Adam Taiye Mudi¹, Bukar Baba Alhaji², Zaharadeen Haruna Aliyu³

¹Department of Mathematics and Statistics, Kaduna Polytechnic, Kaduna, Nigeria

^{2,3}Department of Mathematical Sciences, Nigerian Defence Academy, Kaduna, Nigeria

adamu110202@gmail.com

Abstract

This paper proposes two modified stratified three-stage estimators, denoted as 3M1 and 3M2, for estimating finite population totals under complex survey designs. The proposed estimators aim to improve efficiency by incorporating stratification and Probability Proportional to size (PPS) that reduce bias and variance. Their performance is assessed using both real-life and simulated data sets. Results demonstrate that 3M1 and 3M2 consistently achieve lower variance, mean squared error, and bias compared to conventional estimators such as Horvitz–Thompson (HT), Hansen–Hurwitz (HH), and Hájek Ratio (HR). Relative efficiency measures confirm substantial gains, with 3M1 performing optimally under Scenario 1 and 3M2 under Scenario 2. The study concludes that the proposed estimators provide reliable alternatives for national surveys and large-scale data collection involving stratified multistage designs with 3M1 preferred when balancing variance with bias is required and 3M2 when variance minimization is the priority.

Keywords: Three-stage sampling; Modified estimator; Stratified cluster sampling; Probability proportional to size; Efficiency.

1. Introduction

Large-scale surveys frequently employ stratified multistage cluster sampling to capture population heterogeneity while maintaining operational feasibility and cost control. Three-stage sampling is particularly useful in demographic and health surveys, where primary sampling units (PSUs), secondary sampling units (SSUs), and ultimate units such as households or individuals are selected sequentially. Classical estimators such as the Horvitz–Thompson (HT) estimator and ratio-based variants provide design -unbiasedness, but their efficiency can deteriorate in stratified three-stage settings, especially under probability-proportional-to-size (PPS) selection and heterogeneous clusters. This has led to continuing interest in alternative estimators that balance unbiasedness, variance efficiency, and computational tractability.

The foundational literature on survey sampling, particularly Cochran's Sampling Techniques (1977), remains central to the study of multistage designs, stratification, and the properties of

unbiased estimators. Arnab (2017) extended this by formalizing design-based and model-assisted perspectives on multistage estimation, offering unified treatments of bias, mean squared error (MSE), and efficiency under stratified schemes. Wolter's Introduction to Variance Estimation (2007) further provided methodological tools for replication-based and Taylor-linearized variance estimation, which continue to validate the simulation of classical estimators. Several literatures contributions have sought to address the need for three stage sampling estimation. Nafiu, Oshungade, and Adewara (2012) developed alternative estimation procedures under three-stage sampling designs, demonstrating the possibility of bias reduction relative to HT-type estimators and other classical estimators. Enang and Onyishi (2016) studied variance estimation specifically under three-stage SRSWOR designs, offering insights into design-consistent estimators of sampling variability. More recently, Mudi and Alhaji (2024) introduced modified two-stage cluster sampling estimators, showing improvements in efficiency; their results motivate extensions to stratified three-stage settings.

At the same time, methodological advances in small-area and auxiliary-assisted estimation provide useful parallels. Rao and Molina (2015) highlighted the role of auxiliary information in improving efficiency under complex survey designs, while Lee, Lee, and Shin (2016) showed the potential of composite estimators in stratified two-stage contexts. Although these approaches have not been systematically extended to stratified three-stage estimators, they point toward promising directions for balancing design-unbiasedness and MSE reduction.

Based on these, the present paper introduces two alternative stratified three-stage estimators (3M1 and 3M2) for finite population totals. These estimators integrate ideas from modified and composite estimation into the three-stage case, accommodating PPS selection in the first two stages and SRSWOR at the final stage. Analytical properties including bias, variance, and MSE are derived, and the estimators' relative performance is examined through both Monte Carlo simulations and an application to the Nigeria Demographic and Health Survey DHS 2018 data.

The main aim of this study is to improve efficiency in stratified three-stage cluster sampling for finite population totals. Specifically, the objectives are :to **modify three-stage cluster estimators** by incorporating stratification and PPS at the first and second stages, and SRSWOR at the third stage, using equal probability of selection (3M1) and unequal probability of selection (3M2) methods, to **derive expressions** for the bias, variance, and mean squared error (MSE) of the proposed estimators and to **empirically compare** the performance of the proposed estimators with classical alternatives using both real and simulated data.

The remainder of this paper is organized as follows: Section 2 outlines the survey design framework and reviews classical estimators, while Section 3 introduces the proposed estimators together with their theorems and proofs of statistical properties. Section 4 presents a comparative evaluation of the proposed estimators with conventional estimators commonly used under stratified three stage cluster sampling while Section 5 presents empirical and simulated results for the proposed estimators. Section 6 discusses the results while Section 7 concludes.

2. Methodology

2.1 Survey Design

Consider a finite population U of size N , partitioned into L strata, such that stratum h contains N_h primary sampling units (PSUs), with $\sum_{h=1}^L N_h = N$. Within each stratum h , a sample of n_h PSUs is selected using probability proportional to size (PPS).

Within each selected i^{th} PSU of stratum h , there are M_{hi} secondary sampling units (SSUs). From these, a subsample of m_{hi} SSUs is selected, again using PPS.

Within each selected j^{th} SSU of i^{th} PSU in stratum h , there are M_{hij} tertiary sampling units (TSUs). From these, a final subsample of m_{hij} TSUs is selected using **simple random sampling without replacement (SRSWOR)**. The population total of interest is

$$Y = \sum_{h=1}^L \sum_{i=1}^{N_h} \sum_{j=1}^{M_{hi}} \sum_{\ell=1}^{M_{hij}} y_{hij\ell}, \quad (2.1)$$

where $y_{hij\ell}$ is the value of the study variable for the ℓ^{th} TSU in j^{th} SSU of i^{th} PSU in stratum h . The **inclusion probability** of a TSU under this design is

$$\pi_{hij\ell} = \pi_{hi} \cdot \pi_{hij|i} \cdot \pi_{hij\ell|ij} \quad (2.2)$$

where π_{hi} is the first-stage inclusion probability of i^{th} PSU in stratum h , $\pi_{hij|i}$ is the conditional second-stage probability of selecting j^{th} SSU within the i^{th} PSU and $\pi_{hij\ell|ij}$ is the conditional third-stage probability of selecting ℓ^{th} TSU

2.2 Classical Estimators

Many well-known design-based estimators are employed in multistage sampling among which are:

2.2.1 Horvitz–Thompson (HT) Estimator

$$\hat{Y}_{HT} = \sum_{h=1}^L \sum_{i \in s_h} \sum_{j \in s_{hi}} \sum_{\ell \in s_{hij}} \frac{y_{hij\ell}}{\pi_{hij\ell}}, \quad (2.3)$$

where s_h, s_{hi}, s_{hij} denote the sampled PSUs, SSUs, and TSUs, respectively.

2.2.2 Hansen–Hurwitz (HH) Estimator

In cases where PPS with replacement is used, the HH estimator is given by:

$$\hat{Y}_{HH} = \frac{1}{n} \sum_{r=1}^n \frac{y_r}{p_r}, \quad (2.4)$$

where p_r is the selection probability of the r^{th} draw.

2.2.3 Hájek Ratio (HR) Estimator

The HR estimator improves stability by normalizing with estimated totals of auxiliary measures:

$$\hat{Y}_{HR} = \frac{\sum_{h=1}^L \sum_{i \in s_h} \sum_{j \in s_{hi}} \sum_{\ell \in s_{hij}} \frac{y_{hij\ell}}{\pi_{hij\ell}}}{\sum_{h=1}^L \sum_{i \in s_h} \sum_{j \in s_{hi}} \sum_{\ell \in s_{hij}} \frac{1}{\pi_{hij\ell}}} \cdot N. \quad (2.5)$$

3. Proposed Modified Three-Stage Estimators

To address inefficiency and instability of classical estimators under stratified three-stage designs, two modified estimators are proposed:

(a) Modified Estimator 1 (3M1: Equal Probability at Third Stage)

When equal probability SRSWOR is used at the third stage, the estimator is defined as:

$$\hat{Y}_{3M1} = \sum_{h=1}^L \sum_{i \in s_h} \sum_{j \in s_{hi}} \sum_{\ell \in s_{hij}} W_{hij\ell} y_{hij\ell}, \quad (3.1)$$

where $W_{hij\ell} = \frac{1}{\pi_{hij\ell}}$ is the design weight., and $\pi_{hij\ell} = \pi_{hi} \cdot \pi_{hij/i} \cdot \pi_{hij\ell/ij}$, (the inclusion probabilities at the first, second and third stage) with

$$\pi_{hi} = \frac{n_h M_{hi}}{M_h} \quad (3.2)$$

$$\pi_{hij/i} = \frac{m_h M_{hij}}{M_{hi}} \quad (3.3)$$

$$\pi_{hij\ell/ij} = \frac{m_{hi.}}{M_{hij}} \quad (3.4)$$

Based on the assumption that the clusters are sampled with replacement, an estimator of the variance of \hat{Y}_{3M1} according to Mudi and Alhaji (2024) can be written as

$$V(\hat{Y}_{3M1}) = \sum_{h=1}^L \frac{1}{n_h(n_h - 1)} \sum_{i=1}^{n_h} (y_{hi} - \bar{y}_h)^2 \quad (3.5)$$

$$\text{Where } y_{hi} = \sum_l (n_h w_{hijl}) \hat{Y}_{hijl}, \quad \text{and } \bar{y}_h = \frac{1}{n_h} \sum_i y_{hi}$$

(b) Modified Estimator 2 (3M2: Unequal Probability at Third Stage)

When unequal probability selection is used at the third stage, the estimator is given as:

$$\hat{Y}_{3M2} = \sum_{h=1}^L \sum_{i \in s_h} \sum_{j \in s_{hi}} \sum_{\ell \in s_{hij}} W_{hij\ell} y_{hij\ell}, \quad (3.6)$$

Where π_{hi} and $\pi_{hij/i}$ are the same as in scenario 1 but

$$\pi_{hij\ell} = \frac{m_{hij}}{M_{hij}} \quad (3.7)$$

The variance of \hat{Y}_{3M2} is expressed as:

$$V(\hat{Y}_{3M1}) = \sum_{h=1}^L \frac{1}{n_h(n_h-1)} \sum_{i=1}^{n_h} (y_{hi} - \bar{y}_h)^2 \quad (3.8)$$

where y_{hi} and \bar{y}_h are as defined in equation (3.5) above

Theorem 1 — Design unbiasedness of \hat{Y}_{3M1} and \hat{Y}_{3M2}

Under the sampling scheme in which Stage-1 and Stage-2 use PPS and Stage-3 uses SRSWOR, the estimators

$$\hat{Y}_{3M1} = \sum_{h=1}^L \sum_{i \in s_h} \sum_{j \in s_{hi}} \sum_{\ell \in s_{hij}} w_{hij\ell} y_{hij\ell} \quad (3.9)$$

and

$$\hat{Y}_{3M2} = \sum_{h=1}^L \sum_{i \in s_h} \sum_{j \in s_{hi}} \sum_{\ell \in s_{hij}} w_{hij\ell} y_{hij\ell} \quad (3.10)$$

are design-unbiased estimators of the population total Y . (Both estimators have the same HT form; Scenario 1 vs Scenario 2 differ only in how $\pi_{hij\ell}$ is specified.)

We show the result for \hat{Y}_{3M1} ; the same argument applies to \hat{Y}_{3M2} by replacing the third-stage selection numbers appropriately.

Writing \hat{Y}_{3M1} in HT form we have

$$\hat{Y}_{3M1} = \sum_{h=1}^L \sum_{i \in s_h} \sum_{j \in s_{hi}} \sum_{\ell \in s_{hij}} \frac{y_{hij\ell}}{\pi_{hi} \pi_{hij} \pi_{hij\ell}}. \quad (3.11)$$

Taking the full design expectation $E[\cdot] = E_1[E_2[E_3[\cdot]]]$ and using linearity of expectation and the tower property we have

$$E[\hat{Y}_{3M1}] = E_1 \left[E_2 \left[E_3 \left[\sum_h \sum_{i \in s_h} \sum_{j \in s_{hi}} \sum_{\ell \in s_{hij}} \frac{y_{hij\ell}}{\pi_{hi} \pi_{hij} \pi_{hij\ell}} \right] \right] \right]. \quad (3.12)$$

Evaluating the third-stage expectation E_3 by fixing a sampled j^{th} SSU in sampled i^{th} PSU in stratum h . under SRSWOR of m_{hij} TSUs from M_{hij} TSUs, gives

$$E_3 \left[\sum_{\ell \in s_{hij}} \frac{y_{hij\ell}}{\pi_{hij\ell}} \right] = \sum_{\ell=1}^{M_{hij}} y_{hij\ell}. \quad (3.13)$$

Therefore the inner expectation gives the full SSU total as

$$E_3 \left[\sum_{\ell \in s_{hj}} \frac{y_{hij\ell}}{\pi_{hij\ell}} \right] = Y_{hij} \quad (3.14)$$

where $Y_{hij} = \sum_{\ell=1}^{M_{hij}} y_{hij\ell}$.

Next evaluating the **second-stage** expectation E_2 . Under PPSWR for SSUs within the sampled i^{th} PSU,

$$E_2 \left[\sum_{j \in s_{hi}} \frac{Y_{hij}}{\pi_{hij}} \right] = \sum_{j=1}^{M_{hi}} Y_{hij} = Y_{hi}, \quad (3.15)$$

the total in i^{th} PSU.

Finally evaluating the first-stage expectation E_1 . Under PPSWR r PSUs:

$$E_1 \left[\sum_{i \in s_h} \frac{Y_{hi}}{\pi_{hi}} \right] = \sum_{i=1}^{N_h} Y_{hi}. \quad (3.16)$$

Summing over strata we obtain

$$E[\hat{Y}_{3M1}] = \sum_{h=1}^L \sum_{i=1}^{N_h} \sum_{j=1}^{M_{hi}} \sum_{\ell=1}^{M_{hij}} y_{hij\ell} = Y. \quad (3.17)$$

Hence \hat{Y}_{3M1} is design-unbiased. The same sequence of stagewise expectations holds for \hat{Y}_{3M2} (only the third-stage $\pi_{hij\ell}$ differ), so \hat{Y}_{3M2} is also unbiased.

Theorem 2 — Design-unbiased variance estimation for \hat{Y}_{3M1} and \hat{Y}_{3M2}

$\hat{V}(\hat{Y}_{3M1})$ is a design-unbiased estimator of the design variance of \hat{Y}_{3M1} , that is,

$$E[\hat{V}(\hat{Y}_{3M1})] = V(\hat{Y}_{3M1}).$$

Let us consider a fixed stratum h . The total estimator for stratum h is defined as:

$$\hat{Y}_h = \sum_{i=1}^{n_h} y_{hi}. \quad (3.18)$$

The variance of the stratum total estimator under the sampling design is:

$$Var(\hat{Y}_h) = \sum_{i=1}^{n_h} Var(y_{hi}) + \sum_{i < i'} 2Cov(y_{hi}, y_{hi'}). \quad (3.19)$$

However, assuming that the PSUs are selected independently with replacement, the covariance terms vanish:

$$Cov(y_{hi}, y_{hi'}) = 0, \quad i \neq i'. \quad (3.20)$$

Therefore,

$$\text{Var}(\hat{Y}_h) = n_h \text{Var}(y_{hi}). \quad (3.21)$$

Now, consider the sample variance of the PSU totals in stratum h :

$$\hat{V}_h = \frac{1}{n_h(n_h - 1)} \sum_{i=1}^{n_h} (y_{hi} - \bar{y}_h)^2 \quad (3.22)$$

It is well known that the sample variance \hat{V}_h is an unbiased estimator of the variance of the sample mean \bar{y}_h :

$$E(\hat{V}_h) = \text{Var}(\bar{y}_h). \quad (3.23)$$

Thus, multiplying both sides of (3.23) by n_h^2 yields:

$$E(n_h^2 \hat{V}_h) = n_h^2 \text{Var}(\bar{y}_h) = \text{Var}(\hat{Y}_h). \quad (3.24)$$

Since $\hat{Y}_h = n_h \bar{y}_h$, we can also rewrite:

$$E(\hat{V}_h) = \frac{1}{n_h^2} \text{Var}(\hat{Y}_h) = \frac{1}{n_h} \text{Var}(y_{hi}). \quad (3.25)$$

From (3.21),

$$\text{Var}(y_{hi}) = \frac{1}{n_h} \text{Var}(\hat{Y}_h). \quad (3.26)$$

Plugging (3.26) into (3.25) gives:

$$E(\hat{V}_h) = \frac{1}{n_h} \left(\frac{1}{n_h} \text{Var}(\hat{Y}_h) \right) = \frac{1}{n_h^2} \text{Var}(\hat{Y}_h). \quad (3.27)$$

Since \hat{V}_h estimates the variance of $\hat{Y}_h = n_h \bar{y}_h$, scaling by n_h^2 recovers the correct design scale. Therefore, the aggregate variance estimator is:

$$\hat{V}(\hat{Y}_{3M1}) = \sum_{h=1}^L \hat{V}_h. \quad (3.28)$$

This is an unbiased estimator of the variance of the total estimator:

$$E[\hat{V}(\hat{Y}_{3M1})] = \sum_{h=1}^L E(\hat{V}_h) = \sum_{h=1}^L \text{Var}(\hat{Y}_h) = V(\hat{Y}_{3M1}). \quad (3.29)$$

Similarly,

$$E[\hat{V}(\hat{Y}_{3M2})] = \sum_{h=1}^L E(\hat{V}_h) = \sum_{h=1}^L \text{Var}(\hat{Y}_h) = V(\hat{Y}_{3M2}). \quad (3.30)$$

4. Efficiency Comparison of the Proposed Estimators

This section presents a comparative evaluation of the proposed Modified Three stage estimators 3M1 and 3M2 with conventional estimators commonly used under stratified three stage cluster

sampling. The analysis aims to quantify relative performance in terms of variance, bias, and overall mean square error (MSE).

4.1 Relative Efficiency (RE) Formulation

The **Relative Efficiency (RE)** of an estimator \hat{Y}_{prop} with respect to a comparison estimator \hat{Y}_{comp} is defined as

$$RE(\hat{Y}_{prop}, \hat{Y}_{comp}) = \frac{\text{Var}_p(\hat{Y}_{prop})}{\text{Var}_p(\hat{Y}_{comp})}. \quad (4.1)$$

Values of $RE < 1$ indicate that the proposed estimator is **more efficient** than the comparator, whereas $RE > 1$ implies **less efficiency**.

5. Results

This section presents empirical and simulated results for the proposed estimators. Two scenarios are considered: (i) **real data application** using the DHS Nigeria 2018 extract data, and (ii) **Monte Carlo simulation** studies to assess estimator performance under controlled settings. The analysis was carried using R statistical package to evaluate the performance of the proposed three-stage estimators (3M1 and 3M2) relative to the classical Horvitz–Thompson (HT), Hansen–Hurwitz (HH) and Hájek Ratio (HR). The metrics of evaluation included bias, variance, mean squared error (MSE), and relative efficiency (RE) as presented in the tables below

Table 1: Bias, Variance, RMSE, and MSE of Estimators (DHS and Simulated Data)

Estimator	Dataset	Estimate (\hat{Y})	Bias	Variance ($\times 10^{13}$)	MSE ($\times 10^{13}$)
3M1	DHS	149,401,109	-598,891	2.65	2.69
3M2	DHS	152,269,942	+2,269,942	2.41	2.93
HT	DHS	150,560,098	+560,098	5.00	5.03
HH	DHS	150,139,116	+139,116	15.10	15.10
HR	DHS	145,678,791	-4,321,209	12.20	14.07
3M1	Simulated	5,005	0	120.0	120.0
3M2	Simulated	5,007	+2	125.0	129.0
HT	Simulated	5,000	-5	120.0	145.0
HH	Simulated	5,005	0	140.0	140.0
HR	Simulated	5,002	-3	130.0	139.0

Table 2: Relative Efficiency of Estimators (Benchmark: 3M1 and 3M2)

Estimator	Dataset	RE(3M1=1)	RE(3M2=1)
3M1	DHS	1.00	0.91
3M2	DHS	1.10	1.00
HT	DHS	0.53	0.48
HH	DHS	0.18	0.16
HR	DHS	0.22	0.20
3M1	Simulated	1.00	1.04
3M2	Simulated	0.96	1.00
HT	Simulated	1.00	1.04
HH	Simulated	0.86	0.89
HR	Simulated	0.92	0.96

6. Discussion of Results

The results presented in Tables 1 and 2 represent empirical and simulation results. They demonstrate the advantages of the proposed estimators 3M1 and 3M2 relative to the classical Horvitz-Thompson (HT), Hansen-Hurwitz (HH), and Hajek Ratio (HR) estimators. The discussion is based on how the estimators perform in terms of bias, variance mean square error (MSE) and relative efficiency, using both real data from DHS 2018 and simulated data.

The results presented in Table 1 shows that 3M1 achieves the lowest MSE (2.69×10^{13}) in the DHS application and the same pattern is repeated in the simulation study (120.0). These results suggest that combining stratification with PPS at the first two stages and equal probability SRSWOR at the third stage effectively reduces the overall sampling error.

In terms of variance, For the DHS data only, 3M1 performs well as it attains the smallest variance as compared to the classical estimators considered. 3M2 shows a slightly lower variance (2.41×10^{13}) than 3M1 (2.65×10^{13}).

The bias results further support the reliability of the proposed estimators 3M1 and 3M2 as both exhibit very small bias in the simulated data. This shows the consistency with the theoretical design-based properties established earlier in this paper. In the DHS data, proposed estimators show very small or negligible bias with $-598,891$ for 3M1 and $+2,269,942$ for 3M2.

In table 2, the relative efficiency (RE) values confirm further that using 3M1 and 3M2 as benchmarks, all classical estimators show relative efficient values well below unity ($RE < 1$) indicating that the proposed estimators are more efficient than the classical ones. While 3M2 demonstrates marginal gain in efficiency over 3M1 in certain instances, these gains are offset by increased bias.

7. Conclusions

From the DHS and the simulated data, the proposed estimators (3M1 and 3M2) improve upon classical alternatives by balancing unbiasedness with efficiency. Among them, 3M1 should be preferred when balancing variance with bias is required and 3M2 when variance minimization is the priority.

References

- Arnab, R. (2017). *Survey Sampling Theory and Applications*. Academic Press.
- Cochran, W. G. (1977). *Sampling Techniques* (3rd ed.). John Wiley & Sons.
- Demographic and Health Surveys (DHS). (2018). *Nigeria Demographic and Health Survey 2018*. National Population Commission (NPC) and ICF.
- Enang, E. E., & Onyishi, I. C. (2016). Variance estimation in three-stage sampling designs with simple random sampling without replacement at each stage. *Nigerian Journal of Statistics*, 32(2), 45–60.
- Mudi, T. A., & Alhaji, B. B. (2024). Modified two-stage cluster sampling estimators for finite population total. *Journal of the Royal Statistical Society of Nigeria*, 1(2), 74–89.
- Nafiu, A., Oshungade, I. O., & Adewara, A. A. (2012). Alternative estimation procedures for three-stage sampling designs. *African Journal of Applied Statistics*, 1(2), 33–48.
- Rao, J. N. K., & Molina, I. (2015). *Small Area Estimation* (2nd ed.). Wiley.
- Wolter, K. M. (2007). *Introduction to Variance Estimation* (2nd ed.). Springer.