*Mathematics*

*Research article*

# Efficient estimators of finite population variance using raw moments under two- and three-stage cluster sampling schemes

**Mohsin Abbas**[1]**, Muhammad Ahmed Shehzad**[1]**, Hasnain Iftikhar**[2,3,*]**, Paulo Canas Rodrigues**[4]**, Abdulmajeed Atiah Alharbi**[5]**, and Jeza Allohibi**[5]

[1] Department of Statistics, Bahauddin Zakariya University, Multan, Pakistan

[2] Department of Statistics, University of Peshawar, Peshawar 25120, Pakistan

[3] Department of Statistics, Quaid-i-Azam University, Islamabad 45320, Pakistan

[4] Department of Statistics, Federal University of Bahia, Salvador 40170-110, Brazil

[5] Department of Mathematics, College of Science, Taibah University, Madinah 42353, Saudi Arabia

* **Correspondence:** Email: hasnain@stat.qau.edu.pk.

**Abstract:** In this study, we proposed novel estimators for finite population variance based on the raw moments of the study and auxiliary variables. Specifically, we developed both biased and unbiased estimators of variance using the raw moments of the study variable alone, as well as biased and unbiased difference-type estimators that incorporate the raw moments of a single auxiliary variable. These estimators were evaluated under two-stage cluster sampling (2SCS) and three-stage cluster sampling (3SCS) schemes. Their performance, with and without auxiliary information, was assessed using mean squared error (MSE), absolute bias (AB), and relative efficiency (RE) criteria. Results from two real populations showed that AB decreases and RE improves with increasing sample size. Notably, under 3SCS, the unbiased difference estimator, $\hat{S}^2_{Y,DU}$, achieved the highest efficiency ($RE_3 = 527.69$), closely followed by the biased difference estimator, $\hat{S}^2_{Y,DB}$ ($RE_4 = 527.26$). Both estimators substantially outperformed conventional variance estimators without auxiliary information (baseline $RE = 100$). These findings demonstrate that incorporating auxiliary variables significantly enhances estimation accuracy, offering a practical and robust approach for variance estimation in complex survey designs.

**Keywords:** multi-stage cluster sampling; variance estimation; raw moments, auxiliary information; difference estimator; biased and unbiased estimator; relative efficiency; absolute bias, survey sampling; real-world dataset
**Mathematics Subject Classification:** 62D05, 62G05, 62H12

# 1. Introduction

Accurate estimation of the finite population variance is a cornerstone of survey sampling theory, particularly in the context of complex survey designs such as the 2SCS and 3SCS. These multistage strategies are widely employed in demographic, agricultural, and health surveys due to their operational feasibility and cost-effectiveness. However, the hierarchical structure and intra-cluster correlations inherent in such designs pose significant challenges in achieving efficient and unbiased estimators. Classical approaches, such as those by [1, 2], often rely on restrictive assumptions of cluster homogeneity, which may not hold in real-world applications, potentially leading to inefficiencies and biases.

## 1.1. Classical and design-based approaches

Early contributions to variance estimation focused mainly on simple random sampling (SRS) and cluster sampling frameworks. Traditional estimators developed by [1], extended by [2], and further refined by [3, 4], provided foundational tools but were limited in efficiency under complex multistage structures. Extensions to the 2SCS included variance estimators by [5], which provided exact confidence intervals under equal cluster sizes, and robust approaches by [6, 7], which accounted for within-cluster correlation and imputed data. While these works advanced the field, their reliance on assumptions such as homogeneity or specific model forms has limited their practical applicability.

## 1.2. Auxiliary information and estimator efficiency

The integration of auxiliary information has proven particularly effective in enhancing estimator performance. Auxiliary variables, known for the entire population or significant subpopulations, tend to exhibit a strong correlation with the study variable. Ratio, product, and regression-type estimators using auxiliary data have demonstrated notable reductions in sampling error [2, 8–11]. Within variance estimation, ratio-and regression-based extensions were proposed by [12–15], followed by refinements in [16–20]. More recent developments have considered measurement errors and rank-based auxiliary variables [21–26]. Under SRS, [21, 27, 28] introduced improved and hybrid estimators that significantly outperformed conventional approaches. These contributions collectively establish auxiliary information as a powerful tool to improve estimator efficiency.

## 1.3. Recent developments in complex sampling designs

Despite extensive progress under SRS and stratified designs, fewer contributions directly address variance estimation under complex multistage sampling. Notable exceptions include the robust error variance estimators of [6, 7] for the 2SCS. However, there remains a gap in exploring estimators rooted in raw moment theory and difference-type formulations, particularly when auxiliary information is incorporated. Given the prevalence of multistage sampling in large-scale demographic and health surveys, this gap limits both theoretical advances and practical applications in fields where complex designs are the norm.

## 1.4. Contribution of the present study

To address this gap, the present study proposes a flexible class of biased and unbiased estimators

for finite population variance under the 2SCS and 3SCS designs. The methodology is grounded in raw moment theory, employing both simple and difference-type estimators, with and without auxiliary variables. By systematically incorporating auxiliary data, our approach enhances estimation precision across diverse sampling configurations. We rigorously analyze the bias, MSE, and RE of the proposed estimators and validate their performance using two real-life datasets through extensive simulations.

This framework contributes both theoretically and practically. It generalizes traditional estimators to multistage sampling, provides efficient alternatives when auxiliary information is available, and demonstrates substantial gains in accuracy relative to conventional methods. The results confirm that auxiliary-variable-based estimators consistently outperform their traditional counterparts, highlighting their practical relevance for applications in health, economic, and social surveys where multistage sampling is standard.

The remainder of this paper is organized as follows. Section 2 details the estimation procedures for the finite population mean and its variance under the 2SCS and 3SCS models. Section 3 derives key covariances that inform the structure of the proposed estimators. In Section 4, we present a difference estimator that incorporates auxiliary information. Section 5 outlines computational techniques for variance estimation. Section 6 introduces and evaluates both biased and unbiased estimators. Section 7 presents results from simulation studies and empirical applications, with detailed discussions and study limitations, and highlights potential avenues and methodologies. Finally, Section 8 concludes the study with key insights and directions for future research.

## 2. Estimation of the finite population mean under two-stage and three-stage cluster sampling

This section develops unbiased estimators for the finite population mean $\mu_Y$ under the 2SCS and 3SCS schemes. Further, the variances of these estimators and unbiased estimators of their variances are also derived in this section, which will later be utilized in the variance estimation procedures.

### 2.1. Two-stage cluster sampling

Consider a population consisting of $N$ primary sampling units (PSUs), each containing $N_i$ secondary sampling units (SSUs). Let $Y_{i,j}$ denote the value of the study variable for the $j$th SSU in the $i$th PSU. The finite population mean $\mu_Y$ is defined as:

$$\mu_Y = \frac{1}{N\overline{N}} \sum_{i=1}^{N} N_i \mu_{Y,i}, \tag{2.1}$$

where $\overline{N} = \sum_{i=1}^{N} N_i / N$ is the average cluster size, and $\mu_{Y,i} = \sum_{j=1}^{N_i} Y_{i,j} / N_i$ is the PSU-specific mean.

A sample of $n$ PSUs is selected using SRS without replacement. From each chosen PSU, a sample of $n_i$ SSUs is then drawn accordingly. The 2SCS estimator for $\mu_Y$ is defined as:

$$\overline{Y}_{2S} = \frac{1}{n\overline{N}} \sum_{i=1}^{n} N_i \bar{y}_i, \tag{2.2}$$

where $\bar{y}_i = \sum_{j=1}^{n_i} Y_{i,j} / n_i$ is the sample mean of SSUs within the $i$th PSU. Further elaboration is available in [29] and the references therein.

**Lemma 2.1.** $\overline{Y}_{2S}$ *is an unbiased estimator of* $\mu_Y$.

*Proof.* Let $E_1$ and $E_2$ denote expectations over the first and second sampling stages, respectively. Then,

$$E(\overline{Y}_{2S}) = E_1 \left[ \frac{1}{n\overline{N}} \sum_{i=1}^{n} N_i E_2(\bar{y}_i) \right]$$

$$= E_1 \left[ \frac{1}{n\overline{N}} \sum_{i=1}^{n} N_i \mu_{Y,i} \right]$$

$$= \frac{1}{N\overline{N}} \sum_{i=1}^{N} N_i \mu_{Y,i} = \mu_Y.$$

The result follows from properties of SRS. See also [24, 30] and the references cited therein for more details. □

**Lemma 2.2.** *The variance of* $\overline{Y}_{2S}$ *is given by:*

$$V(\overline{Y}_{2S}) = \frac{\lambda \sigma_{Y,b,2}^2}{n} + \frac{1}{nN\overline{N}^2} \sum_{i=1}^{N} \frac{\lambda_i N_i^2 \sigma_{Y,i,2}^2}{n_i}, \tag{2.3}$$

*where* $\lambda = 1 - n/N$, $\lambda_i = 1 - n_i/N_i$, *and*

$$\sigma_{Y,b,2}^2 = \frac{1}{N-1} \sum_{i=1}^{N} (N_i \mu_{Y,i} - \overline{N}\mu_Y)^2,$$

$$\sigma_{Y,i,2}^2 = \frac{1}{N_i - 1} \sum_{j=1}^{N_i} (Y_{i,j} - \mu_{Y,i})^2.$$

*Proof.* By the law of total variance,

$$V(\overline{Y}_{2S}) = V_1 \left[ E_2(\overline{Y}_{2S}) \right] + E_1 \left[ V_2(\overline{Y}_{2S}) \right].$$

Using the properties of expectations,

$$V_1 \left[ E_2(\overline{Y}_{2S}) \right] = \frac{\lambda \sigma_{Y,b,2}^2}{n\overline{N}^2},$$

$$E_1 \left[ V_2(\overline{Y}_{2S}) \right] = \frac{1}{nN\overline{N}^2} \sum_{i=1}^{N} \frac{\lambda_i N_i^2 \sigma_{Y,i,2}^2}{n_i}.$$

Combining the two components yields the stated result. □

**Lemma 2.3.** *An unbiased estimator of* $V(\overline{Y}_{2S})$ *is:*

$$\widehat{V}(\overline{Y}_{2S}) = \frac{\lambda \hat{\sigma}_{Y,b,2}^2}{n} + \frac{1}{nN\overline{N}^2} \sum_{i=1}^{N} \frac{\lambda_i N_i^2 \hat{\sigma}_{Y,i,2}^2}{n_i - 1}, \tag{2.4}$$

*where*

$$\hat{\sigma}^2_{Y,b,2} = \frac{1}{n-1} \sum_{i=1}^{n} (N_i \bar{y}_i - \overline{N}\, \overline{Y}_{2S})^2, \tag{2.5}$$

$$\hat{\sigma}^2_{Y,i,2} = \frac{1}{n_i-1} \sum_{j=1}^{n_i} (Y_{i,j} - \bar{y}_i)^2. \tag{2.6}$$

*Proof.* From Eq (2.5), we can write

$$\hat{\sigma}^2_{Y,2b} = \frac{n}{n-1} \left[ \frac{1}{n} \sum_{i=1}^{n} (N_i \bar{Y}_i)^2 - (\overline{N}\, \overline{Y}_{2S})^2 \right]. \tag{2.7}$$

Take mathematical expectations of the terms on the right side of Eq (2.7) to get

$$
\begin{aligned}
E\left[ \frac{1}{n} \sum_{i=1}^{n} (N_i \bar{Y}_i)^2 \right] &= E_1 \left[ \frac{1}{n} \sum_{i=1}^{n} E_2 (N_i \bar{Y}_i)^2 \right] \\
&= E_1 \left[ \frac{1}{n} \sum_{i=1}^{n} \left( V_2(N_i \bar{Y}_i) + \{E_2(N_i \bar{Y}_i)\}^2 \right) \right] \\
&= E_1 \left[ \frac{1}{n} \sum_{i=1}^{n} \left\{ \frac{\lambda_i M_i^2 \sigma^2_{Y,i,2}}{n_i} + M_i^2 \mu^2_{Y,i} \right\} \right] \\
&= \frac{1}{N} \sum_{i=1}^{N} \frac{\lambda_i M_i^2 \sigma^2_{Y,i,2}}{n_i} + \frac{1}{N} \sum_{i=1}^{N} M_i^2 \mu^2_{Y,i}
\end{aligned}
$$

and

$$
\begin{aligned}
E(\overline{N}\, \overline{Y}_{2S})^2 &= V(\overline{N}\, \overline{Y}_{2S}) + \left( E(\overline{N}\, \overline{Y}_{2S}) \right)^2 \\
&= \frac{\lambda \sigma^2_{Y,b,2}}{n} + \frac{1}{nN} \sum_{i=1}^{N} \frac{\lambda_i M_i^2 \sigma^2_{Y,i,2}}{n_i} + (\overline{N} \mu_Y)^2.
\end{aligned}
$$

Thus, $\hat{\sigma}^2_{Y,b,2}$ is a biased estimator of $\sigma^2_{Y,b,2}$. Moreover, from Eq (2.6), we can write:

$$\hat{\sigma}^2_{Y,2i} = \frac{n_i}{n_i-1} \left[ \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}^2 - (\bar{Y}_i)^2 \right]. \tag{2.8}$$

Take mathematical expectations of terms on the right side of Eq (2.8) to get

$$
\begin{aligned}
E_2\left[ \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}^2 \right] &= \left[ V_2(\frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}) + \left( E_2(\frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}) \right)^2 \right] \\
&= \frac{\sigma^2_{Y,i,2}}{n_i} + \frac{1}{N_i} \sum_{j=1}^{N_i} Y_{ij}^2
\end{aligned}
$$

and

$$
\begin{aligned}
E_2(\bar{Y}_i)^2 &= V_2(\bar{Y}_i) + [E_2(\bar{Y}_i)]^2 \\
&= \frac{\lambda_i \sigma_{Y,i,2}^2}{n_i} + (\mu_{Y,i})^2.
\end{aligned}
$$

Thus, $\hat{\sigma}_{Y,i,2}^2$ is an unbiased estimator of $\sigma_{Y,i}^2, 2$. Now, by Eq (2.7), we already have the expectation of $\sum_{i=1}^{n}(N_i\bar{Y}_i)^2/n$, and by Eq (2.8) we know the expectation of within-cluster variance $\hat{\sigma}_{Y,i,2}^2$. Substituting these results into Eq (2.4), we get

$$
E\left(\widehat{V}(\bar{Y}_{2S})\right) = \frac{n}{n-1}\left\{E\left[\frac{1}{n}\sum_{i=1}^{n}(N_i\bar{Y}_i)^2\right] - E\left[(\overline{N}\bar{Y}_{2S})^2\right]\right\}. \tag{2.9}
$$

Now, substituting the results of Eqs (2.7) and (2.8) into the above and grouping like terms, we obtain

$$
E\left(\widehat{V}(\bar{Y}_{2S})\right) = \frac{n}{n-1}\left[\left(1-\frac{1}{n}\right)\frac{1}{N}\sum_{i=1}^{N}\frac{\lambda_i N_i^2 \sigma_{Y,i,2}^2}{n_i} + \left(\frac{1}{N}\sum_{i=1}^{N}(N_i\mu_{Y,i})^2 - (\overline{N}\mu_Y)^2\right) - \frac{\lambda\sigma_{Y,b,2}^2}{n}\right]. \tag{2.10}
$$

Recall that the between-cluster variance is defined as

$$
\sigma_{Y,b,2}^2 = \frac{1}{N-1}\sum_{i=1}^{N}\left(N_i\mu_{Y,i} - \overline{N}\mu_Y\right)^2. \tag{2.11}
$$

Now, consider the cross-term appearing in Eq (2.10),

$$
\frac{1}{N}\sum_{i=1}^{N}(N_i\mu_{Y,i})^2 - (\overline{N}\mu_Y)^2. \tag{2.12}
$$

Expanding this using the definition in Eq (2.11), we have

$$
\begin{aligned}
\frac{1}{N}\sum_{i=1}^{N}(N_i\mu_{Y,i})^2 - (\overline{N}\mu_Y)^2 &= \frac{1}{N}\sum_{i=1}^{N}\left[(N_i\mu_{Y,i} - \overline{N}\mu_Y)^2 + 2(N_i\mu_{Y,i} - \overline{N}\mu_Y)(\overline{N}\mu_Y) + (\overline{N}\mu_Y)^2\right] - (\overline{N}\mu_Y)^2 \\
&= \frac{1}{N}\sum_{i=1}^{N}(N_i\mu_{Y,i} - \overline{N}\mu_Y)^2 + \frac{2\overline{N}\mu_Y}{N}\sum_{i=1}^{N}(N_i\mu_{Y,i} - \overline{N}\mu_Y) + \frac{1}{N}\sum_{i=1}^{N}(\overline{N}\mu_Y)^2 - (\overline{N}\mu_Y)^2.
\end{aligned}
$$

Since $\sum_{i=1}^{N}(N_i\mu_{Y,i} - \overline{N}\mu_Y) = 0$, the middle term vanishes, and the last two terms cancel each other. Thus,

$$
\frac{1}{N}\sum_{i=1}^{N}(N_i\mu_{Y,i})^2 - (\overline{N}\mu_Y)^2 = \frac{1}{N}\sum_{i=1}^{N}(N_i\mu_{Y,i} - \overline{N}\mu_Y)^2. \tag{2.13}
$$

Using Eq (2.12), this simplifies to

$$
\frac{1}{N}\sum_{i=1}^{N}(N_i\mu_{Y,i})^2 - (\overline{N}\mu_Y)^2 = \frac{N-1}{N}\sigma_{Y,b,2}^2. \tag{2.14}
$$

Substituting Eq (2.14) into the Eq (2.10), and combining with the within-cluster component, we obtain

$$E\left(\widehat{V}(\bar{Y}_{2S})\right) = \frac{\lambda \sigma_{Y,b,2}^2}{n} + \frac{1}{nN\overline{N}^2} \sum_{i=1}^{N} \frac{\lambda_i N_i^2 \sigma_{Y,i,2}^2}{n_i}, \tag{2.15}$$

which establishes the desired result.

$$E\left(\widehat{V}(\bar{Y}_{2S})\right) = V(\bar{Y}_{2S}) \tag{2.16}$$

which completes the proof. The proof of this lemma is also provided in [29, 31], where a similar derivation is presented. □

### 2.2. The 3SCS scheme

The 3SCS methodology improves upon the 2SCS approach by adding an extra stage to the sampling process. In this extended 3SCS structure, the overall population $U$ consists of $N$ PSUs, each of which contains $N_i$ SSUs, and each SSU is further subdivided into $N_{ij}$ tertiary sampling units (TSUs). The characteristic value for the $k$th TSU within the $j$th SSU of the $i$th PSU is denoted as $Y_{ijk}$, where $i$ ranges from 1 to $N$, $j$ ranges from 1 to $N_i$, and $k$ ranges from 1 to $N_{ij}$. Within this framework, the population mean, $\mu_Y$, is defined as

$$\mu_Y = \frac{1}{N\overline{T}} \sum_{i=1}^{N} \sum_{j=1}^{N_i} N_{ij} \mu_{Y,ij}, \tag{2.17}$$

where $\overline{T} = \sum_{i=1}^{N} \sum_{j=1}^{N_i} N_{ij}/N$ represents the average cluster size, while $\mu_{Y,ij} = \sum_{k=1}^{N_{ij}} Y_{ij,k}/N_{ij}$ denotes the mean calculated from the $j$th SSU of the $i$th PSU.

In the 3SCS framework, the population mean $\mu_Y$ is estimated by deriving an unbiased estimator through a structured sampling procedure. Initially, a sample of PSUs, denoted as $n$, is selected. Subsequently, a sample of SSUs, $n_i$, is drawn from each PSU, followed by the selection of TSUs, $n_{ij}$, within each SSU. It is important to note that SRS without replacement is employed at each stage of the sampling under 3SCS scheme. The estimator for $\mu_Y$ is then given by:

$$\bar{Y}_{3S} = \frac{1}{n\overline{T}} \sum_{i=1}^{n} \frac{N_i}{n_i} \sum_{j=1}^{n_i} N_{ij} \bar{y}_{ij} \quad , \tag{2.18}$$

where $\bar{y}_{ij} = \sum_{k=1}^{n_{ij}} Y_{ij,k}/n_{ij}$.

In the forthcoming lemmas, we explore the mathematical attributes of $\bar{Y}_{3S}$.

**Lemma 2.4.** $\bar{Y}_{3S}$ *is an unbiased estimators of* $\mu_Y$.

*Proof.* Here, 1–3, with $E$ and $V$, refer to the expectation and variance at the first, second, and third stages of sampling, respectively. By assigning these, we obtain:

$$E(\bar{Y}_{3S}) = E_1 E_2 \left[ \frac{1}{n\overline{T}} \sum_{i=1}^{n} \frac{N_i}{n_i} \sum_{j=1}^{n_i} N_{ij} E_3(\bar{y}_{ij}) \right]$$

$$
\begin{aligned}
&= E_1 E_2 \left[ \frac{1}{n\overline{T}} \sum_{i=1}^{n} \frac{N_i}{n_i} \sum_{j=1}^{n_i} N_{ij}\mu_{Y,ij} \right] \\
&= E_1 \left[ \frac{1}{n\overline{T}} \sum_{i=1}^{n} \sum_{j=1}^{N_i} N_{ij}\mu_{Y,ij} \right] \\
&= \frac{1}{N\overline{T}} \sum_{i=1}^{N} \sum_{j=1}^{N_i} N_{ij}\mu_{Y,ij} = \mu_Y,
\end{aligned}
\tag{2.19}
$$

which completes the proof. $\square$

**Lemma 2.5.** *The variances of $\overline{Y}_{3S}$ is given by:*

$$
V(\overline{Y}_{3S}) = \frac{\lambda\sigma_{Y,b,3}^2}{n\overline{T}^2} + \frac{1}{nN\overline{T}^2} \sum_{i=1}^{N} \frac{\lambda_i N_i^2 \sigma_{Y,i,3}^2}{n_i} + \frac{1}{nN\overline{T}^2} \sum_{i=1}^{N} \frac{n_i}{n_i} \sum_{j=1}^{N_i} \frac{N_{ij}^2 \lambda_{ij} \sigma_{Y,ij}^2}{n_{ij}},
\tag{2.20}
$$

*where*

$$
\sigma_{Y,b,3}^2 = \frac{1}{N-1} \sum_{i=1}^{N} (N_i\mu_{Y,i} - \overline{T}\mu_Y)^2, \quad \sigma_{Y,i,3}^2 = \frac{1}{N_i-1} \sum_{j=1}^{N_i} \left( N_{ij}\mu_{Y,ij} - \mu_{Y,i} \right)^2,
$$

$$
\sigma_{Y,ij}^2 = \frac{1}{N_{ij}-1} \sum_{k=1}^{N_{ij}} \left( Y_{ij,k} - \mu_{Y,ij} \right)^2, \quad \mu_{Y,i} = \frac{1}{N_i} \sum_{j=1}^{N_i} N_{ij}\mu_{Y,ij}, \quad and \quad \lambda_{ij} = \left( 1 - \frac{n_{ij}}{N_{ij}} \right).
$$

In the 3SCS, between-cluster variation denoted by $\sigma_{Y,b,3}^2$, known as the first-stage variation, refers to the differences among the PSUs, capturing the diversity of clusters. Between-sub-cluster variation denoted by $\sigma_{Y,i,3}^2$, known as the second-stage variation, reflects the differences among SSUs within each selected PSU, highlighting variation at the sub-cluster level. Within-sub-cluster variation denoted by $\sigma_{Y,ij}^2$, known as the third-stage variation, involves the variability within TSUs selected from each SSU, capturing the final level of within-cluster diversity.

*Proof.* As it is known that

$$
V(\overline{Y}_{3S}) = E_1 E_2 V_3 [\overline{Y}_{3S}] + E_1 V_2 E_3 [\overline{Y}_{3S}] + V_1 E_2 E_3 [\overline{Y}_{3S}],
\tag{2.21}
$$

from Eq (2.19), we can write

$$
\begin{aligned}
V_1 E_2 E_3 [\overline{Y}_{3S}] &= V_1 E_2 \left[ \frac{1}{n\overline{T}} \sum_{i=1}^{n} \frac{N_i}{n_i} \sum_{j=1}^{n_i} N_{ij}\mu_{Y,ij} \right] \\
&= V_1 \left( \frac{1}{n\overline{T}} \sum_{i=1}^{n} N_i\mu_{Y,i} \right) = \frac{\lambda\sigma_{Y,b,3}^2}{n\overline{T}^2}
\end{aligned}
\tag{2.22}
$$

and

$$
E_1 V_2 E_3 [\overline{Y}_{3S}] = E_1 V_2 \left[ \frac{1}{n\overline{T}} \sum_{i=1}^{n} N_i \bar{y}_i \right]
$$

$$= E_1 \left[ \frac{1}{n_1^2 \overline{T}^2} \sum_{i=1}^{n} N_i^2 \, V_2(\bar{y}_i) \right]$$

$$= \frac{1}{nN\overline{T}^2} \sum_{i=1}^{N} \frac{\lambda_i N_i^2 \sigma_{Y,i,3}^2}{n_i}, \tag{2.23}$$

where $\bar{y}_i = \sum_{j=1}^{n_i} N_{ij} \mu_{Y,ij} / n_i$. Also from Eq (2.19), we have:

$$\begin{aligned} V_3[\overline{Y}_{3S}] &= \frac{1}{n^2 \overline{T}^2} \sum_{i=1}^{n} \frac{N_i^2}{n_i^2} \sum_{j=1}^{n_i} N_{ij}^2 \, V_3(\bar{y}_{ij}) \\ &= \frac{1}{n^2 \overline{T}^2} \sum_{i=1}^{n} \frac{N_i^2}{n_i^2} \sum_{j=1}^{n_i} \frac{N_{ij}^2 \lambda_{ij} \sigma_{Y,ij}^2}{n_{ij}}. \end{aligned} \tag{2.24}$$

We take expectations of Eq (2.24) to get

$$E_1 E_2 V_3[\overline{Y}_{3S}] = \frac{1}{nN\overline{T}^2} \sum_{i=1}^{N} \frac{N_i}{n_i} \sum_{j=1}^{N_i} \frac{N_{ij}^2 \lambda_{ij} \sigma_{Y,ij}^2}{n_{ij}}. \tag{2.25}$$

Eq (2.20) is derived by substituting Eqs (2.22), (2.23), and (2.25) into Eq (2.21), thereby concluding the proof. $\square$

**Lemma 2.6.** *Unbiased estimators of $V(\overline{Y}_{3S})$ are given as follows:*

$$\widehat{V}(\overline{Y}_{3S}) = \frac{\lambda \hat{\sigma}_{Y,b,3}^2}{n\overline{T}^2} + \frac{1}{nN\overline{T}^2} \sum_{i=1}^{n} \frac{\lambda_i N_i^2 \hat{\sigma}_{Y,i,3}^2}{n_i} + \frac{1}{nN\overline{T}^2} \sum_{i=1}^{n} \frac{N_i}{n_i} \sum_{j=1}^{n_i} \frac{N_{ij}^2 \lambda_{ij} \hat{\sigma}_{Y,ij}^2}{n_{ij}}, \tag{2.26}$$

*respectively, where*

$$\hat{\sigma}_{Y,b,3}^2 = \frac{1}{n-1} \sum_{i=1}^{n} \left( N_i \bar{y}_i - \overline{T} \, \overline{Y}_{3S} \right)^2, \tag{2.27}$$

$$\hat{\sigma}_{Y,i,3}^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} \left( N_{ij} \bar{y}_{ij} - \bar{y}_i \right)^2, \tag{2.28}$$

$$\hat{\sigma}_{Y,ij}^2 = \frac{1}{n_{ij} - 1} \sum_{k=1}^{n_{ij}} \left( Y_{ij,k} - \bar{y}_{ij} \right)^2. \tag{2.29}$$

*Proof.* From Eq (2.27), we can write

$$\hat{\sigma}_{Y,b,3}^2 = \frac{n}{n-1} \left[ \frac{1}{n} \sum_{i=1}^{n} (N_i \bar{y}_i)^2 - (\overline{T} \, \overline{Y}_{3S})^2 \right]. \tag{2.30}$$

Consider the expected values of the terms on the right-hand side (RHS) of Eq (2.30) to obtain:

$$E \left[ \frac{1}{n} \sum_{i=1}^{n} (N_i \bar{y}_i)^2 \right] = E_1 E_2 \left[ \frac{1}{n} \sum_{i=1}^{n} E_3 (N_i \bar{y}_i)^2 \right]$$

$$
\begin{aligned}
&= E_1 E_2 \left[ \frac{1}{n} \sum_{i=1}^{n} \left\{ V_3(N_i \bar{y}_i) + (E_3(N_i \bar{y}_i))^2 \right\} \right] \\
&= E_1 E_2 \left[ \frac{1}{n} \sum_{i=1}^{n} \left\{ \frac{N_i^2}{n_{i,2}^2} \sum_{j=1}^{n_i} N_{ij}^2 V_3(\bar{y}_{ij}) + (N_i \bar{y}_i)^2 \right\} \right] \\
&= E_1 \left[ \frac{1}{n} \sum_{i=1}^{n} \left\{ \frac{N_i}{n_i} \sum_{j=1}^{N_i} \frac{\lambda_{ij} N_{ij}^2 \sigma_{Y,ij}^2}{n_{ij}} + \frac{\lambda_i N_i^2 \sigma_{Y,i,3}^2}{n_i} + (N_i \mu_{Y,i})^2 \right\} \right] \\
&= \frac{1}{N} \sum_{i=1}^{N} \frac{N_i}{n_i} \sum_{j=1}^{N_i} \frac{\lambda_{ij} N_{ij}^2 \sigma_{Y,ij}^2}{n_{ij}} + \frac{1}{N} \sum_{i=1}^{N} \frac{\lambda_i N_i^2 \sigma_{Y,i,3}^2}{n_i} + \frac{1}{N} \sum_{i=1}^{N} (N_i \mu_{Y,i})^2 \quad (2.31)
\end{aligned}
$$

and

$$
\begin{aligned}
E[\overline{T}\ \overline{Y}_{3S}]^2 &= V[\overline{T}\ \overline{Y}_{3S}] + \left[ E(\overline{T}\ \overline{Y}_{3S}) \right]^2 \\
&= \frac{\lambda \sigma_{Y,b,3}^2}{n} + \frac{1}{nN} \sum_{i=1}^{N} \frac{N_i}{n_i} \sum_{j=1}^{N_i} \frac{\lambda_{ij} N_{ij}^2 \sigma_{Y,ij}^2}{n_{ij}} + \frac{1}{nN} \sum_{i=1}^{N} \frac{\lambda_i N_i^2 \sigma_{Y,i,3}^2}{n_i} + (\overline{T}\ \mu_Y)^2. \quad (2.32)
\end{aligned}
$$

By combining Eqs (2.31) and (2.32), and subsequently taking the expected values of Eq (2.27), we obtain:

$$
\begin{aligned}
E(\hat{\sigma}_{Y,b,3}^2) &= \frac{n}{n-1} \left[ \frac{1}{N} \sum_{i=1}^{N} \frac{N_i}{n_i} \sum_{j=1}^{N_i} \frac{\lambda_{ij} N_{ij}^2 \sigma_{Y,ij}^2}{n_{ij}} + \frac{1}{N} \sum_{i=1}^{N} \frac{\lambda_i N_i^2 \sigma_{Y,i,3}^2}{n_i} + \frac{1}{N} \sum_{i=1}^{N} (N_i \mu_{Y,i})^2 \right. \\
&\quad \left. - \left\{ \frac{\lambda \sigma_{Y,b,3}^2}{n} + \frac{1}{nN} \sum_{i=1}^{N} \frac{N_i}{n_i} \sum_{j=1}^{N_i} \frac{\lambda_{ij} N_{ij}^2 \sigma_{Y,ij}^2}{n_{ij}} + \frac{1}{nN} \sum_{i=1}^{N} \frac{\lambda_i N_i^2 \sigma_{Y,i,3}^2}{n_i} + (\overline{T}\ \mu_Y)^2 \right\} \right]. \quad (2.33)
\end{aligned}
$$

Now grouping like terms, we can rewrite the above as

$$
\begin{aligned}
E(\hat{\sigma}_{Y,b,3}^2) &= \frac{n}{n-1} \left[ \left( 1 - \frac{1}{n} \right) \frac{1}{N} \sum_{i=1}^{N} \frac{N_i}{n_i} \sum_{j=1}^{N_i} \frac{\lambda_{ij} N_{ij}^2 \sigma_{Y,ij}^2}{n_{ij}} \right. \\
&\quad \left. + \left( 1 - \frac{1}{n} \right) \frac{1}{N} \sum_{i=1}^{N} \frac{\lambda_i N_i^2 \sigma_{Y,i,3}^2}{n_i} + \left\{ \frac{1}{N} \sum_{i=1}^{N} (N_i \mu_{Y,i})^2 - (\overline{T}\ \mu_Y)^2 \right\} - \frac{\lambda \sigma_{Y,b,3}^2}{n} \right].
\end{aligned}
$$

Thus, after simplification and reorganization as given in Lemma 2.3, we obtain the final biased expected value:

$$
E(\hat{\sigma}_{Y,b,3}^2) = \sigma_{Y,b,3}^2 + \frac{1}{N} \sum_{i=1}^{N} \frac{\lambda_i N_i^2 \sigma_{Y,i,3}^2}{n_i} + \frac{1}{N} \sum_{i=1}^{N} \frac{N_i}{n_i} \sum_{j=1}^{N_i} \frac{\lambda_{ij} N_{ij}^2 \sigma_{Y,ij}^2}{n_{ij}}. \quad (2.34)
$$

This demonstrates that $\hat{\sigma}_{Y,b,3}^2$ is a biased estimator of $\sigma_{Y,b,3}^2$. For more details, see [30, 32]. Similarly, we can derive from Eq (2.28):

$$
\hat{\sigma}_{Y,i,3}^2 = \frac{n_i}{(n_i - 1)} \left[ \frac{1}{n_i} \sum_{j=1}^{n_i} (N_{ij} \bar{y}_{ij})^2 - (\bar{y}_i)^2 \right]. \quad (2.35)
$$

Now take the expected values of the terms on the RHS of Eq (2.35) to obtain:

$$
\begin{aligned}
E_2\left[\frac{1}{n_i}\sum_{j=1}^{n_i}(N_{ij}\bar{y}_{ij})^2\right] &= E_2\left[\frac{1}{n_i}\sum_{j=1}^{n_i}E_3(N_{ij}\bar{y}_{ij})^2\right] \\
&= E_2\left[\frac{1}{n_i}\sum_{j=1}^{n_i}\left\{V_3(N_{ij}\bar{y}_{ij})+\left(E_3(N_{ij}\bar{y}_{ij})\right)^2\right\}\right] \\
&= \frac{1}{N_i}\sum_{j=1}^{N_i}\frac{\lambda_{ij}N_{ij}^2\sigma_{Y,ij}^2}{n_{ij}}+\frac{1}{N_i}\sum_{j=1}^{N_i}\left(N_{ij}\mu_{Y,ij}\right)^2
\end{aligned}
\tag{2.36}
$$

and

$$
\begin{aligned}
E_2(\bar{y}_i)^2 &= E_2E_3(\bar{y}_i)^2 \\
&= E_2\left[V_3(\bar{y}_i)+(E_3(\bar{y}_{i,\mathrm{SRS}}))^2\right] \\
&= E_2\left[\frac{1}{n_{i,2}^2}\sum_{j=1}^{n_i}\frac{\lambda_{ij}N_{ij}^2\sigma_{Y,ij}^2}{n_{ij}}+\left\{\frac{1}{n_i}\sum_{j=1}^{n_i}N_{ij}\mu_{Y,ij}\right\}^2\right] \\
&= \frac{1}{n_iN_i}\sum_{j=1}^{N_i}\frac{\lambda_{ij}N_{ij}^2\sigma_{Y,ij}^2}{n_{ij}}+E_2\left\{\frac{1}{n_i}\sum_{j=1}^{n_i}N_{ij}\mu_{Y,ij}\right\}^2 \\
&= \frac{1}{n_iN_i}\sum_{j=1}^{N_i}\frac{\lambda_{ij}N_{ij}^2\sigma_{Y,ij}^2}{n_{ij}}+V_2\left\{\frac{1}{n_i}\sum_{j=1}^{n_i}N_{ij}\mu_{Y,ij}\right\}+\left\{E_2(\frac{1}{n_i}\sum_{j=1}^{n_i}N_{ij}\mu_{Y,ij})\right\}^2 \\
&= \frac{\lambda_i\sigma_{Y,i,3}^2}{n_i}+\frac{1}{n_iN_i}\sum_{j=1}^{N_i}\frac{\lambda_{ij}N_{ij}^2\sigma_{Y,ij}^2}{n_{ij}}+(\mu_{Y,i})^2.
\end{aligned}
\tag{2.37}
$$

By combining Eqs (2.36) and (2.37), and subsequently taking the expected values of Eq (2.35), we obtain:

$$
E_2\left(\hat{\sigma}_{Y,i,3}^2\right) = \sigma_{Y,i,3}^2+\frac{1}{N_i}\sum_{j=1}^{N_i}\frac{\lambda_{ij}N_{ij}^2\sigma_{Y,ij}^2}{n_{ij}}.
\tag{2.38}
$$

This indicates that $\hat{\sigma}_{Y,i,3}^2$ serves as a biased estimator of $\sigma_{Y,i,3}^2$. Similarly, we can derive from Eq (2.29):

$$
\hat{\sigma}_{Y,ij}^2 = \frac{n_{ij}}{n_{ij}-1}\left[\frac{1}{n_{ij}}\sum_{k=1}^{n_{ij}}Y_{ij,k}^2-(\bar{y}_{ij})^2\right].
\tag{2.39}
$$

Consider the expectation of the terms on the RHS of Eq (2.39), which demonstrate that $\hat{\sigma}_{Y,ij}^2$ serves as an unbiased estimator of $\sigma_{Y,ij}^2$, i.e.,

$$
E_3(\hat{\sigma}_{Y,ij}^2) = \sigma_{Y,ij}^2.
\tag{2.40}
$$

By following the same lines provided in Lemma 2.3, consider the expectation of Eq (2.26), and apply the results from Eqs (2.34), (2.38), and (2.40) to show that:

$$
E\left(\widehat{V}(\overline{Y}_{3S})\right) = V(\overline{Y}_{3S}),
\tag{2.41}
$$

which completes the proof. □

In the previous section, we developed unbiased estimators for the population mean $\mu_Y$, along with their variances and unbiased estimators of variance under the 2SCS and 3SCS schemes. It is important to note that the estimation for the auxiliary variable $X$ follows exactly the same procedure as for $Y$. In the next section, we extend this framework to derive the covariance between the mean estimators and its unbiased estimator. Estimating this covariance is essential for constructing difference and regression-type estimators, which are subsequently used to enhance the efficiency of variance estimation.

## 3. Covariance computation and estimation under the 2SCS and 3SCS schemes

In this section, we employ the previously mentioned 2SCS/3SCS schemes to derive accurate mathematical formulas for the covariances and unbiased estimator of covariance of the mean estimators, which are utilized in the following sections. The covariances $C(\bar{Y}_{2S}, \bar{X}_{2S})$ and $C(\bar{Y}_{3S}, \bar{X}_{3S})$, along with their unbiased estimators $\widehat{C}(\bar{Y}_{2S}, \bar{X}_{2S})$ and $\widehat{C}(\bar{Y}_{3S}, \bar{X}_{3S})$, serve as key components in improving the efficiency of difference and regression-type estimators by incorporating the correlation between the study and auxiliary variables. It is important to note that SRS without replacement is applied at each stage of sampling under the 2SCS/3SCS schemes, and it is assumed that information on both the study and auxiliary variables is available and observations within and between clusters are independent with homogeneous variance. These assumptions ensure the theoretical properties of the estimators, such as unbiasedness and MSE, hold.

### 3.1. The 2SCS scheme

In a finite population $\Omega$, we consider $X$ as the auxiliary variable and $Y$ as the study variable of interest. To estimate $(\mu_Y, \mu_X)$ within the 2SCS framework, let $(\bar{Y}_{2S}, \bar{X}_{2S})$ denote the mean estimators obtained from the variables $(Y, X)$. It is worth noting that the mean estimators for $X$ can be calculated in a manner similar to those for $Y$.

**Lemma 3.1.** *Under the 2SCS scheme, the covariance between $\bar{Y}_{2S}$ and $\bar{X}_{2S}$ and its unbiased estimator are expressed as:*

$$C(\bar{Y}_{2S}, \bar{X}_{2S}) = \frac{\lambda \sigma_{XY,b,2}}{n\bar{N}^2} + \frac{1}{nN\bar{N}^2} \sum_{i=1}^{N} \frac{\lambda_i N_i^2 \sigma_{XY,i,2}}{n_i}, \tag{3.1}$$

$$\widehat{C}(\bar{Y}_{2S}, \bar{X}_{2S}) = \frac{\lambda \hat{\sigma}_{XY,b,2}}{n\bar{N}^2} + \frac{1}{nN\bar{N}^2} \sum_{i=1}^{n} \frac{\lambda_i M_i^2 \hat{\sigma}_{XY,i,2}}{n_i}, \tag{3.2}$$

*where*

$$\sigma_{XY,b,2} = \frac{1}{N-1} \sum_{i=1}^{N} \left( (N_i \mu_{Y,i} - \bar{N} \mu_Y)(N_i \mu_{X,i} - \bar{N} \mu_X) \right),$$

$$\sigma_{XY,i,2} = \frac{1}{N_i - 1} \sum_{j=1}^{N_i} \left( (Y_{i,j} - \mu_{Y,i})(X_{i,j} - \mu_{X,i}) \right),$$

$$\hat{\sigma}_{XY,b,2} = \frac{1}{n-1} \sum_{i=1}^{n} \left( (N_i \bar{Y}_i - \bar{N} \bar{Y}_{2S})(N_i \bar{X}_i - \bar{N} \bar{X}_{2S}) \right),$$

$$\hat{\sigma}_{XY,i,2} \;=\; \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (Y_{i,j} - \bar{Y}_i)(X_{i,j} - \bar{X}_i).$$

Here, $\sigma_{XY,b,2}$, $\hat{\sigma}_{XY,b,2}$, and $\sigma_{XY,i,2}$, $\hat{\sigma}_{XY,i,2}$ *have their usual meanings.*

*Proof.* Here, 1 and 2, with $C$ and $E$, refer to the covariance and expectation at the first and second stages of sampling, respectively. By assigning these, we obtain the covariance between $\overline{Y}_{2S}$ and $\overline{X}_{2S}$:

$$C(\overline{Y}_{2S}, \overline{X}_{2S}) \;=\; C_1[E_2(\overline{Y}_{2S}, \overline{X}_{2S})] + E_1[C_2(\overline{Y}_{2S}, \overline{X}_{2S})].$$

It can be shown that $E_2(\overline{Y}_{2S}) = \sum_{i=1}^{n} N_i \mu_{Y,i}/(n\overline{N})$. Based on this result, we have

$$C_1\left[E_2(\overline{Y}_{2S}, \overline{X}_{2S})\right] \;=\; C_1\left(\frac{1}{n\overline{N}}\sum_{i=1}^{n} N_i\mu_{Y,i}, \frac{1}{n\overline{N}}\sum_{i=1}^{n} N_i\mu_{X,i}\right) \;=\; \frac{\lambda\sigma_{XY,b,2}}{n\overline{N}^2}, \tag{3.3}$$

$$E_1\left[C_2(\overline{Y}_{2S}, \overline{X}_{2S})\right] \;=\; E_1\left[\frac{1}{n^2\overline{N}^2}\sum_{i=1}^{n} N_i^2\, C_2(\mu_{Y,i}, \mu_{X,i})\right],$$

$$\;=\; \frac{1}{nN\overline{N}^2}\sum_{i=1}^{N} \frac{\lambda_i N_i^2 \sigma_{XY,i,2}}{n_i}. \tag{3.4}$$

Add Eqs (3.3) and (3.4), and this completes the proof. $\qquad\square$

### 3.2. The 3SCS scheme

In a finite population $\Omega$, we consider $X$ as the auxiliary and $Y$ as the study variable of interest. To estimate $(\mu_Y, \mu_X)$ within the context of the 3SCS, let $(\overline{Y}_{3S}, \overline{X}_{3S})$ denote the respective mean estimators obtained from the variables $(Y, X)$.

**Lemma 3.2.** *In the 3SCS framework, the covariance between $\overline{Y}_{3S}$ and $\overline{X}_{3S}$ including its unbiased estimator are given as follows:*

$$C(\overline{Y}_{3S}, \overline{X}_{3S}) \;=\; \frac{\lambda\sigma_{XY,b,3}}{n\overline{T}^2} + \frac{1}{nN\overline{T}^2}\sum_{i=1}^{N}\frac{\lambda_i N_i^2 \sigma_{XY,i,3}}{n_i} + \frac{1}{nN\overline{T}^2}\sum_{i=1}^{N}\frac{N_i}{n_i}\sum_{j=1}^{N_i}\frac{\lambda_{ij} N_{ij}^2 \sigma_{XY,ij,3}}{n_{ij}}, \tag{3.5}$$

$$\widehat{C}(\bar{Y}_{3S}, \bar{X}_{3S}) \;=\; \frac{\lambda\hat{\sigma}_{XY,b,3}}{n\overline{T}^2} + \frac{1}{nN\overline{T}^2}\sum_{i=1}^{n}\frac{\lambda_i M_i^2 \hat{\sigma}_{XY,i,3}}{n_i} + \frac{1}{nN\overline{T}^2}\sum_{i=1}^{n}\frac{N_i}{n_i}\sum_{j=1}^{n_i}\frac{\lambda_{ij} T_{ij}^2 \hat{\sigma}_{XY,ij,3}}{n_{ij}}, \tag{3.6}$$

*where*

$$\sigma_{XY,b,3} \;=\; \frac{1}{N-1}\sum_{i=1}^{N}\Big((N_i\mu_{Y,i} - \overline{T}\mu_Y)(N_i\mu_{X,i,3} - \overline{T}\mu_X)\Big),$$

$$\hat{\sigma}_{XY,3b} \;=\; \frac{1}{n-1}\sum_{i=1}^{n}\Big(N_i\bar{Y}_i - \overline{N}\bar{Y}_{3S})(N_i\bar{X}_i - \overline{N}\bar{X}_{3S})\Big),$$

$$\sigma_{XY,i,3} \;=\; \frac{1}{N_i-1}\sum_{j=1}^{N_i}\Big((N_{ij}\mu_{Y,ij} - \mu_{Y,i})(N_{ij}\mu_{X,ij} - \mu_{X,i,3})\Big),$$

$$\hat{\sigma}_{XY,3i} = \frac{1}{m_i - 1} \sum_{j=1}^{m_i} \left( (N_{ij}\bar{Y}_{ij} - \bar{Y}_i)(N_{ij}\bar{X}_{ij} - \bar{X}_i) \right),$$

$$\sigma_{XY,ij,3} = \frac{1}{N_{ij} - 1} \sum_{k=1}^{N_{ij}} \left( (Y_{ij,k} - \mu_{Y,ij})(X_{ij,k} - \mu_{X,ij}) \right),$$

$$\hat{\sigma}_{XY,3ij} = \frac{1}{n_{ij} - 1} \sum_{k=1}^{n_{ij}} \left( (Y_{ij,k} - \bar{Y}_{ij})(X_{ij,k} - \bar{X}_{ij}) \right),$$

and $\lambda_{ij} = (1 - n_{ij}/N_{ij})$, where $\sigma_{XY,b,3}$, $\sigma_{XY,i,3}$, and $\sigma_{XY,ij,3}$ have their usual meanings. For more details, the reader may see [33] and the references cited therein.

*Proof.* Here, 1–3, with $C$ and $E$, refer to the covariance and expectation at the first, second, and third stages of sampling, respectively. By assigning these, we obtain the covariance between $\bar{Y}_{3S}$ and $\bar{X}_{3S}$:

$$\begin{aligned}
C(\bar{Y}_{3S}, \bar{X}_{3S}) &= C_1 E_2 E_3 \left[ \bar{Y}_{3S}, \bar{X}_{3S} \right] + E_1 C_2 E_3 \left[ \bar{Y}_{3S}, \bar{X}_{3S} \right] \\
&\quad + E_1 E_2 C_3 \left[ \bar{Y}_{3S}, \bar{X}_{3S} \right].
\end{aligned} \tag{3.7}$$

It can be demonstrated, as shown in Eq (2.19), that $E_3(\bar{Y}_{3S}) = \sum_{i=1}^{n}(N_i/n_i) \sum_{j=1}^{n_i} N_{ij}\mu_{Y,ij}/n\bar{T}$. Consequently, we obtain:

$$\begin{aligned}
C_1 E_2 E_3 \left[ \bar{Y}_{3S}, \bar{X}_{3S} \right] &= C_1 E_2 \left[ \frac{1}{n\bar{T}} \sum_{i=1}^{n} \frac{N_i}{n_i} \sum_{j=1}^{n_i} N_{ij}\mu_{Y,ij}, \frac{1}{n\bar{T}} \sum_{i=1}^{n} \frac{N_i}{n_i} \sum_{j=1}^{n_i} N_{ij}\mu_{X,ij} \right] \\
&= C_1 \left[ \frac{1}{n\bar{T}} \sum_{i=1}^{n} N_i\mu_{Y,i}, \frac{1}{n\bar{T}} \sum_{i=1}^{n} N_i\mu_{X,i,3} \right] = \frac{\lambda\sigma_{XY,b,3}}{n\bar{T}^2}.
\end{aligned} \tag{3.8}$$

$$\begin{aligned}
E_1 C_2 E_3 \left[ \bar{Y}_{3S}, \bar{X}_{3S} \right] &= E_1 C_2 \left[ \frac{1}{n\bar{T}} \sum_{i=1}^{n} \frac{N_i}{n_i} \sum_{j=1}^{n_i} N_{ij}\mu_{Y,ij}, \frac{1}{n\bar{T}} \sum_{i=1}^{n} \frac{N_i}{n_i} \sum_{j=1}^{n_i} N_{ij}\mu_{X,ij} \right] \\
&= E_1 \left[ \frac{1}{n^2\bar{T}^2} \sum_{i=1}^{n} N_i^2 \, C_2 \left( \frac{1}{n_i} \sum_{j=1}^{n_i} N_{ij}\mu_{Y,ij}, \frac{1}{n_i} \sum_{j=1}^{n_i} N_{ij}\mu_{X,ij} \right) \right] \\
&= E_1 \left[ \frac{1}{n^2\bar{T}^2} \sum_{i=1}^{n} \frac{\lambda_i N_i^2 \sigma_{XY,i,3}}{n_i} \right] \\
&= \frac{1}{nN\bar{T}^2} \sum_{i=1}^{N} \frac{\lambda_i N_i^2 \sigma_{XY,i,3}}{n_i}.
\end{aligned} \tag{3.9}$$

$$\begin{aligned}
E_1 E_2 C_3 \left[ \bar{Y}_{3S}, \bar{X}_{3S} \right] &= E_1 E_2 \left[ \frac{1}{n^2\bar{T}^2} \sum_{i=1}^{n} \frac{N_i^2}{n_i^2} \sum_{j=1}^{n_i} N_{ij}^2 \, C_3(\bar{y}_{ij}, \bar{x}_{ij}) \right] \\
&= E_1 E_2 \left[ \frac{1}{n^2\bar{T}^2} \sum_{i=1}^{n} \frac{N_i^2}{n_i^2} \sum_{j=1}^{n_i} \frac{\lambda_{ij} N_{ij}^2 \sigma_{XY,ij,3}}{n_{ij}} \right] \\
&= E_1 \left[ \frac{1}{n^2\bar{T}^2} \sum_{i=1}^{n} \frac{N_i}{n_i} \sum_{j=1}^{N_i} \frac{\lambda_{ij} N_{ij}^2 \sigma_{XY,ij,3}}{n_{ij}} \right]
\end{aligned}$$

$$= \frac{1}{nN\overline{T}^2} \sum_{i=1}^{N} \frac{N_i}{n_i} \sum_{j=1}^{N_i} \frac{\lambda_{ij} N_{ij}^2 \sigma_{XY,ij,3}}{n_{ij}}. \tag{3.10}$$

Add Eqs (3.8)–(3.10), and this completes the proof. For more details, see [30]. The proofs of $\widehat{C}(\bar{Y}_{2S}, \bar{X}_{2S})$, and $\widehat{C}(\bar{Y}_{3S}, \bar{X}_{3S})$ can be seen in the study of [33], where a similar derivation is present. $\square$

In the previous section, we derived the covariances and their unbiased estimators under the 2SCS and 3SCS schemes. In the next section, we focus on deriving the difference estimator. The relative error terms, which will be used to express the bias and MSE of the proposed estimators, are also introduced here. This difference estimator will subsequently be employed in the variance estimation procedures to enhance the efficiency of the estimators.

## 4. Difference estimator

This section examines the unbiased difference estimator for the population mean and second raw moment using study data and an auxiliary variable. These estimators were used to construct biased and unbiased variance estimators within an S-sampling framework. The biases and MSEs of the proposed estimators for $S_Y^2$ are evaluated through the following relative error terms:

$$\xi_0 = \frac{\overline{Y}_S - \mu_Y}{\mu_Y} \quad \text{and} \quad \xi_1 = \frac{\overline{X}_S - \mu_X}{\mu_X},$$

$$\xi_0' = \frac{\overline{Y}'_S - \mu'_Y}{\mu'_Y} \quad \text{and} \quad \xi_1' = \frac{\overline{X}'_S - \mu'_X}{\mu'_X},$$

such that $E(\xi_0) = E(\xi_1) = E(\xi_0') = E(\xi_1') = 0$. Let us denote

$$V_{rs} = E(\xi_0^r \xi_1^s) = E\left[\left(\frac{\overline{Y}_S - \mu_Y}{\mu_Y}\right)^r \left(\frac{\overline{X}_S - \mu_X}{\mu_X}\right)^s\right],$$

$$V_{rs}' = E(\xi_0'^r \xi_1'^s) = E\left[\left(\frac{\overline{Y}'_S - \mu'_Y}{\mu'_Y}\right)^r \left(\frac{\overline{X}'_S - \mu'_X}{\mu'_X}\right)^s\right],$$

which gives

$$V_{20} = E(\xi_0^2) = \frac{V(\overline{Y}_S)}{\mu_Y^2}, \qquad\qquad V_{20}' = E(\xi_0'^2) = \frac{V(\overline{Y}'_S)}{\mu_Y'^2},$$

$$V_{02} = E(\xi_1^2) = \frac{V(\overline{X}_S)}{\mu_X^2}, \qquad\qquad V_{02}' = E(\xi_1'^2) = \frac{V(\overline{X}'_S)}{\mu_X'^2},$$

$$V_{11} = E(\xi_0 \xi_1) = \frac{C(\overline{Y}_S, \overline{X}_S)}{\mu_Y \mu_X}, \qquad\qquad V_{11}' = E(\xi_0' \xi_1') = \frac{C(\overline{Y}'_S, \overline{X}'_S)}{\mu_Y' \mu_X'},$$

$$V_{11}^{(*)} = E(\xi_0 \xi_0') = \frac{C(\overline{Y}'_S, \overline{Y}_S)}{\mu_Y \mu_X}.$$

Here, $\overline{Y}_S$ and $\overline{X}_S$ represent the mean estimators of the study and auxiliary variables $Y$ and $X$, respectively, while $\overline{Y}'_S$ and $\overline{X}'_S$ denote the mean estimators of $Y^2$ and $X^2$, respectively, based on an S-sampling scheme, where S is either 2S or 3S.

**Lemma 4.1.** *The difference estimators of $E(Y)$ and $E(Y^2)$ using information on $(Y, X)$ are given by*

$$\overline{Y}_D = \overline{Y}_S + \beta(\mu_X - \overline{X}_S), \tag{4.1}$$

$$\overline{Y}'_D = \overline{Y}'_S + \beta'(\mu'_X - \overline{X}'_S), \tag{4.2}$$

*respectively, where*

$$\beta_S = \frac{C(\overline{Y}_S, \overline{X}_S)}{V(\overline{X}_S)} \quad and \quad \beta'_S = \frac{C(\overline{Y}'_S, \overline{X}'_S)}{V(\overline{X}'_S)}$$

*are known constants.*

*Proof.* $\overline{Y}_D$ and $\overline{Y}'_D$ are the unbiased estimators of $\mu_Y$, and $\mu'_Y$, respectively [2, 34]. The variance of $\overline{Y}_D$ is derived by rewriting Eq (4.1) in terms of $\varepsilon$'s. To calculate the $V(\overline{Y}_D)$, we represent it in relation to $\varepsilon$'s, i.e.,

$$\bar{Y}_D = \mu_Y(1 + \varepsilon_0) - \beta_S\mu_X\varepsilon_1$$

$$\bar{Y}_D - \mu_Y = \mu_Y\varepsilon_0 - \beta_S\mu_X\varepsilon_1. \tag{4.3}$$

To calculate the variance of $\bar{Y}_D$, we square both sides of Eq (4.3) and apply the expectation.

$$V(\bar{Y}_D) = \mu_Y^2\Lambda_{20} + \beta_S^2\mu_X^2\Lambda_{02} - 2\beta_S\mu_X\mu_Y\Lambda_{11}. \tag{4.4}$$

The simplified equation for the variance of $\bar{Y}_D$ may be obtained by substituting $\beta_S$. The variance of $\overline{Y}_D$ is given by

$$V(\overline{Y}_D) = \mu_Y^2 V_{20}\left(1 - \frac{V_{11}^2}{V_{20}V_{02}}\right) \tag{4.5}$$

$$= \mu_Y^2 V_{20}(1 - \varrho^2)$$

where $\varrho = V_{11}/\sqrt{V_{20}V_{02}}$ represents the correlation coefficient between $\overline{Y}_S$ and $\overline{X}_S$ within the context of the S-sampling scheme. On the same lines, one can prove the variance of $\overline{Y}'_D$ given in Eq (4.2). □

The difference estimator can be viewed as a special case of the regression estimator. Specifically, when the regression coefficient $\hat{\beta}_S$ is chosen as the population regression coefficient (or its unbiased estimate), the regression estimator reduces to the difference estimator: $\bar{Y}_D = \bar{Y}_S + \beta(\mu_X - \bar{X}_S)$. This shows that the difference estimator essentially applies a linear adjustment using auxiliary information, achieving variance reduction under the condition that $\hat{\beta}_S$ correctly reflects the relationship between $Y$ and $X$. The value of $\beta_S$ may be obtained through prior research, surveys, or census data. If this value is not available, it may be approximated using a large sample size. In order to determine the estimated value of $\beta_S$, one can substitute $C(\bar{Y}_S, \bar{X}_S)$ and $V(\bar{X}_S)$ with their respective unbiased estimators.

$$\widehat{\beta}_S = \frac{\widehat{C}(\overline{Y}_S, \overline{X}_S)}{\widehat{V}(\overline{X}_S)} \quad and \quad \beta'_S = \frac{\widehat{C}(\overline{Y}'_S, \overline{X}'_S)}{\widehat{V}(\overline{X}'_S)}.$$

The sample covariance estimator, $\widehat{C}(\bar{Y}_S, \bar{X}_S)$, $\widehat{C}(\bar{Y}'_S, \bar{X}'_S)$, and the sample variance estimator, $\widehat{V}(\bar{X}_S)$, $\widehat{V}(\bar{X}'_S)$, are known to exhibit weak consistency in estimating $C(\bar{Y}_S, \bar{X}_S)$, $C(\bar{Y}'_S, \bar{X}'_S)$, and $V(\bar{X}_S)$, $V(\bar{X}'_S)$, respectively, under an SRS scheme. Consequently, for larger sample sizes, $\hat{\beta}_S$ can be regarded as a weakly consistent estimator of $\beta_S$, see [30, 35] for more details. Since the regression-type estimators rely on the adjustment factor $\hat{\beta}$, it is necessary to first estimate the covariance between the study and auxiliary variables. The covariance estimator $\widehat{C}(\bar{Y}_S, \bar{X}_S)$ ensures that $\hat{\beta} = \widehat{C}(\bar{Y}_S, \bar{X}_S)/\widehat{V}(\bar{X}_S)$ can be computed from the sample, thereby allowing auxiliary information to be incorporated. This step is essential because it accounts for the correlation structure between $Y$ and $X$, which leads to improved efficiency compared to estimators that ignore auxiliary variables. The estimator $\bar{Y}_D$, when incorporating the estimated value of $\beta_S$, is commonly referred to as the regression estimator.

$$\bar{Y}_{Reg} = \bar{Y}_S + \hat{\beta}_S \left( \mu_X - \bar{X}_S \right).$$

Similarly, for $Y^2$, the regression estimator can be given as follows:

$$\bar{Y}'_{Reg} = \bar{Y}'_S + \hat{\beta}'_S \left( \mu'_X - \bar{X}'_S \right).$$

It can be demonstrated that the estimators $\bar{Y}_{Reg}$ and $\bar{Y}'_{Reg}$ are biased in estimating $\mu_Y$. Furthermore, as the sample size increases significantly, the following can be observed for $\bar{Y}_{Reg}$:

$$\text{MSE}(\bar{Y}_{Reg}) \approx V(\bar{Y}_D) = \mu_Y^2 \Lambda_{200}(1 - \rho^2). \tag{4.6}$$

Similarly, the variance of $\bar{Y}'_{Reg}$ can be obtained from the above expression. In the next section, we present the expressions for computing the finite population variance under the 2SCS and 3SCS schemes.

## 5. Computation of the finite population variance under the 2SCS and 3SCS schemes

This section provides an overview of the finite population variance, $S_Y^2$, under the 2SCS/3SCS, which will be utilized in the subsequent sections.

### 5.1. The 2SCS scheme

The following expression can be used to obtain the finite population variance under the 2SCS:

$$
\begin{aligned}
S_Y^2 &= \frac{1}{N\overline{N} - 1} \sum_{i=1}^{N} \sum_{j=1}^{N_i} (Y_{ij} - \mu_Y)^2 \\
&= \frac{N\overline{N}}{N\overline{N} - 1} \left[ \frac{1}{N\overline{N}} \sum_{i=1}^{N} \sum_{j=1}^{N_i} Y_{ij}^2 - (\mu_Y)^2 \right] = \eta \left[ \mu'_Y - (\mu_Y^2) \right] \\
&= \eta \left[ E(Y^2) - (E(Y))^2 \right],
\end{aligned} \tag{5.1}
$$

where $\eta = N\overline{N}/(N\overline{N} - 1)$. The aim is to formulate difference estimators for $E(Y)$ and $E(Y^2)$ within the framework of the 2SCS, which will subsequently be employed to derive unbiased estimators of the finite population variance, $S_Y^2$. For more details, refer to [24] and the cited references.

## 5.2. The 3SCS scheme

Under the 3SCS, the expression below calculates the finite population variance:

$$
\begin{aligned}
S_Y^2 &= \frac{1}{N\overline{T}-1} \sum_{i=1}^{N} \sum_{j=1}^{N_i} \sum_{k=1}^{N_{ij}} (Y_{ij,k} - \mu_Y)^2 \\
&= \frac{N\overline{T}}{N\overline{T}-1} \left[ \frac{1}{N\overline{T}} \sum_{i=1}^{N} \sum_{j=1}^{N_i} \sum_{k=1}^{N_{ij}} Y_{ij,k}^2 - \mu_Y^2 \right] = \eta \left[ \mu_Y' - \mu_Y^2 \right] \\
&= \eta \left[ E(Y^2) - (E(Y))^2 \right],
\end{aligned}
\tag{5.2}
$$

where $\eta = N\overline{T}/(N\overline{T}-1)$.

In the next section, we derive the proposed estimators for the finite population variance. These include both biased and unbiased estimators, constructed with and without the use of auxiliary information.

## 6. Proposed estimators

This section presents both biased and unbiased estimators of the finite population variance, $S_Y^2$, developed with and without auxiliary information, under an S-sampling scheme. These estimators are constructed using the first and second raw moments of the study and auxiliary variables to enhance efficiency and accuracy.

### 6.1. First proposed estimator

Following the framework established by [24], we have introduced both biased and unbiased estimators for the variance of a finite population.

$$
\hat{S}_{Y,SB}^2 = \eta \left[ \overline{Y}_S' - (\overline{Y}_S)^2 \right],
\tag{6.1}
$$

$$
\hat{S}_{Y,SU}^2 = \eta \left[ \overline{Y}_S' - (\overline{Y}_S)^2 + \widehat{V}(\overline{Y}_S) \right].
\tag{6.2}
$$

Based solely on information $Y$ under an S-sampling scheme. In this context, $\overline{Y}_S'$ represents the second moment (i.e., the sample mean of squared observations), while $\overline{Y}_S$ denotes the first moment (i.e., the sample mean of the observations) under the S-sampling scheme.

In the following lemmas, we establish the mathematical properties of the proposed estimators $\hat{S}_{Y,SB}^2$, $\hat{S}_{Y,SU}^2$, $\hat{S}_{Y,DB}^2$, and $\hat{S}_{Y,DU}^2$, including their biased and unbiased nature. We also derive the expressions for their variances to assess their efficiency.

**Lemma 6.1.** $\hat{S}_{Y,SB}^2$ is an unbiased estimator of $\hat{S}_Y^2$.

*Proof.* The mathematical expectation of $\hat{S}_{Y,SB}^2$ is

$$
E(\hat{S}_{Y,SB}^2) = \eta \left[ E(\overline{Y}_S') - E(\overline{Y}_S^2) \right]
$$

$$= \eta\left[\mu'_Y - \mu_Y^2 - V(\bar{Y}_S)\right]$$
$$= S_Y^2 - \eta V(\bar{Y}_S),$$

which shows that $\hat{S}_{Y,SB}^2$ underestimates $S_Y^2$; hence, it is a biased estimator of $S_Y^2$. $\qquad\square$

**Lemma 6.2.** *The variance of $\hat{S}_{Y,SB}^2$ under S-sampling scheme is given as follows:*

$$V(\hat{S}_{Y,SB}^2) = \eta^2\left[Var(\bar{Y}'_S) + 4\mu_Y^2 Var(\bar{Y}_S) - 4\mu_Y Cov(\bar{Y}'_S, \bar{Y}_S)\right]$$
$$= \eta^2\left[\mu'^2_Y V'_{20} + 4\mu_Y^4 V_{20} - 4\mu_Y^2\mu'_Y V_{11}^{(*)}\right]. \tag{6.3}$$

To derive the variances of the proposed finite population variance estimators, we use a linearization approach with relative error terms, $\xi_0 = (\bar{Y}_S - \mu_Y)/\mu_Y$ for linear observations and $\xi'_0 = (\bar{Y}'_S - \mu'_Y)/\mu'_Y$ for squared observations. This allows variance decomposition incorporating the variance of linear ($V_{20}$) and squared ($V'_{20}$) observations, as well as their covariance ($V_{11}^{(*)}$), which accounts for their interaction. For difference-type estimators, auxiliary information reduces variance through correlations with study variables, and a finite population correction is applied to account for the sampling design. See Table 1 for more details.

**Table 1.** Variance and covariance terms in Lambda notation used for all proposed estimators.

| Term | Definition | Lambda notation |
|------|------------|-----------------|
| $V(\bar{Y}_S)$ | Variance of sample mean under S | $\mu_Y^2\Lambda_{20}$ |
| $V(\bar{Y}'_S)$ | Variance of the mean of squared observations | $\mu_{Y^2}^2\Lambda_{40}$ |
| $V(\bar{Y}_D)$ | Variance of difference mean | $\mu_Y^2 V_{20}(1-\varrho^2)$ |
| $V(\bar{Y}'_D)$ | Variance of difference mean of squared | $\mu'^2_Y V'_{20}(1-\varrho'^2)$ |
| $C(\bar{Y}'_S, \bar{Y}_S)$ | Covariance between $\bar{Y}'_S$ and $\bar{Y}_S$ | $\mu_Y\mu'_Y V_{11}^{(*)}$ |
| $C(\bar{Y}'_D, \bar{Y}_D)$ | Covariance between $\bar{Y}'_D$ and $\bar{Y}_D$ | $\mu_Y\mu'_Y V_{11}^{(*)}(1-\varrho\varrho')$ |

*Proof.* Express $\hat{S}_{Y,SB}^2$ in terms of relative errors:

$$\hat{S}_{Y,SB}^2 = \eta\left[\mu'_Y(1+\xi'_0) - (\mu_Y(1+\xi_0))^2\right]. \tag{6.4}$$

Expanding the squared term:

$$(\mu_Y(1+\xi_0))^2 = \mu_Y^2(1+2\xi_0+\xi_0^2) \approx \mu_Y^2(1+2\xi_0), \tag{6.5}$$

where the higher-order term $\xi_0^2$ is negligible under the first-order approximation. Substituting back, we obtain the linearized form:

$$\hat{S}_{Y,SB}^2 \approx \eta\left[\mu'_Y(1+\xi'_0) - \mu_Y^2(1+2\xi_0)\right] = \eta\left[(\mu'_Y - \mu_Y^2) + \mu'_Y\xi'_0 - 2\mu_Y^2\xi_0\right]. \tag{6.6}$$

Using the linearized form of the estimator:

$$\hat{S}_{Y,SB}^2 \approx \eta\left[(\mu'_Y - \mu_Y^2) + \mu'_Y\xi'_0 - 2\mu_Y^2\xi_0\right], \tag{6.7}$$

we calculate its variance. By the properties of variance, a constant term does not contribute:

$$V(\hat{S}^2_{Y,SB}) = V\big(\eta(\mu'_Y\xi'_0 - 2\mu^2_Y\xi_0)\big) = \eta^2 V(\mu'_Y\xi'_0 - 2\mu^2_Y\xi_0). \tag{6.8}$$

The term $(\mu'_Y - \mu^2_Y)$ is constant with respect to the sample. Variance measures the variability of a random quantity around its mean, and constants do not vary. Therefore, it does not contribute to the variance. Only the deviations from the mean, captured by the relative errors $\xi_0$ and $\xi'_0$, determine the variability of the estimator. Higher-order terms like $\xi^2_0$ are also ignored under the first-order linearization approximation.

Using the variance formula for a linear combination of random variables:

$$V(aX + bY) = a^2 Var(X) + b^2 V(Y) + 2ab C(X, Y), \tag{6.9}$$

we expand $V(\mu'_Y\xi'_0 - 2\mu^2_Y\xi_0)$ as

$$V(\mu'_Y\xi'_0 - 2\mu^2_Y\xi_0) = (\mu'_Y)^2 V(\xi'_0) + (2\mu^2_Y)^2 V(\xi_0) - 2 \cdot 2\mu^2_Y \cdot \mu'_Y C(\xi'_0, \xi_0). \tag{6.10}$$

Substituting the definitions of the relative error variances and covariance:

$$V'_{20} = V(\xi'_0), \quad V_{20} = V(\xi_0), \quad V^{(*)}_{11} = C(\xi'_0, \xi_0), \tag{6.11}$$

we obtain

$$V(\hat{S}^2_{Y,SB}) = \eta^2\Big[\mu'^2_Y V'_{20} + 4\mu^4_Y V_{20} - 4\mu^2_Y\mu'_Y V^{(*)}_{11}\Big], \tag{6.12}$$

which completes the proof. □

**Lemma 6.3.** $\hat{S}^2_{Y,SU}$ *is an unbiased estimator of* $S^2_Y$.

*Proof.* By taking the expectation of Eq (6.2), we obtain:

$$\begin{aligned} E(\hat{S}^2_{Y,SU}) &= E(\hat{S}^2_{Y,SB}) + \eta\, E(\widehat{V}(\overline{Y}_S)) \\ &= E(\hat{S}^2_{Y,SB}) + \eta\, V(\overline{Y}_S) = S^2_Y. \end{aligned}$$

Hence, $\hat{S}^2_{Y,SU}$ is an unbiased estimator of $\hat{S}^2_Y$ under S-sampling scheme. □

**Lemma 6.4.** *The variance of* $\hat{S}^2_{Y,SU}$ *is given as follows:*

$$V(\hat{S}^2_{Y,SU}) = \eta^2\Big[(\mu'_Y)^2 V'_{20} + 4\mu^4_Y V_{20} - 4\mu^2_Y\mu'_Y V^*_{11}\Big]. \tag{6.13}$$

*Proof.* The second estimator is defined as the unbiased version of the first estimator:

$$\hat{S}^2_{Y,SU} = \hat{S}^2_{Y,SB} + \eta V(\overline{Y}_S) = \hat{S}^2_{Y,SB} + \eta\mu^2_Y V_{20}. \tag{6.14}$$

Since the term $\eta\mu^2_Y V_{20}$ is constant with respect to the sample, it does not contribute to the variance. Therefore,

$$V(\hat{S}^2_{Y,SU}) = V(\hat{S}^2_{Y,SB} + \eta\mu^2_Y V_{20}). \tag{6.15}$$

Substituting the variance of the first estimator, we obtain:

$$V(\hat{S}^2_{Y,SU}) = \eta^2\Big[(\mu'_Y)^2 V'_{20} + 4\mu^4_Y V_{20} - 4\mu^2_Y\mu'_Y V^*_{11}\Big] \approx V(\hat{S}^2_{Y,SB}), \tag{6.16}$$

which completes the proof. □

Remark: Adding a constant to an estimator does not change its variance. Hence, the variance of the second (unbiased) estimator is the same as that of the first estimator.

## 6.2. Second proposed estimator

In line with the methodology outlined by [24], we introduce biased and unbiased estimators for the finite population variance within an S-sampling framework, leveraging data from a single auxiliary variable, $X$.

$$\hat{S}^2_{Y,DB} = \eta\left[\overline{Y}'_D - \overline{Y}^2_D\right], \tag{6.17}$$

$$\hat{S}^2_{Y,DU} = \hat{S}^2_{Y,DB} + \eta V(\overline{Y}_D). \tag{6.18}$$

**Lemma 6.5.** $\hat{S}^2_{Y,DB}$ *is a biased estimator of* $\hat{S}^2_Y$.

*Proof.* The mathematical expectation of $\hat{S}^2_{Y,DB}$ is

$$\begin{aligned}
E(\hat{S}^2_{Y,DB}) &= \eta\left[E(\overline{Y}'_D) - E(\overline{Y}^2_D)\right] \\
&= \eta\left[\mu'_Y - \mu^2_Y - V(\overline{Y}_D)\right] \\
&= S^2_Y - \eta\mu^2_Y V_{20}(1 - \varrho^2),
\end{aligned}$$

which also shows that $\hat{S}^2_{Y,DB}$ underestimates $S^2_Y$, and it is thus a biased estimator of $S^2_Y$. $\qquad\square$

The estimator, $\hat{S}^2_{Y,DB} = \eta\left(\overline{Y}'_D - (\overline{Y}_D)^2\right)$, is a difference-type biased estimator that incorporates auxiliary information. The difference-type sample means are defined as $\overline{Y}_D = \overline{Y}_S - \beta(\overline{X}_S - \mu_X)$ and $\overline{Y}'_D = \overline{Y}'_S - \beta'(\overline{X}'_S - \mu'_X)$, where $\beta$ and $\beta'$ are the regression coefficients between the study variable and auxiliary variable for linear and squared observations, respectively.

**Lemma 6.6.** *The variance of* $\hat{S}^2_{Y,DB}$ *under an S-sampling is given as follows:*

$$V(\hat{S}^2_{Y,DB}) = \eta^2\left[\mu'^2_Y V'_{20}(1 - \varrho'^2) + 4\mu^4_Y V_{20}(1 - \varrho^2) - 4\mu^2_Y\mu'_Y V^{(*)}_{11}(1 - \varrho\varrho')\right]. \tag{6.19}$$

To derive the variance, we first linearize $\overline{Y}_D$ and $\overline{Y}'_D$ using relative error terms $\xi_0, \xi'_0, \xi_1, \xi'_1$, which allows us to center the estimators at zero. Using these error terms, the variances of the adjusted means are reduced due to the auxiliary information, giving $Var(\overline{Y}_D) = \mu^2_Y V_{20}(1 - \varrho^2)$ and $Var(\overline{Y}'_D) = \mu'^2_Y V'_{20}(1 - \varrho'^2)$, where $\varrho$ and $\varrho' = V'_{11}/\sqrt{V'_{20}V'_{02}}$, are the correlations between the study and auxiliary variables. Then the linearized form of the estimator is:

*Proof.*

$$\hat{S}^2_{Y,DB} \approx \eta\left[(\mu'_Y - \mu^2_Y) + \mu'_Y\xi'_0 - 2\mu^2_Y\xi_0\right]. \tag{6.20}$$

For a linear combination of random variables $a\xi'_0 + b\xi_0$:

$$V(a\xi'_0 + b\xi_0) = a^2 V(\xi'_0) + b^2 V(\xi_0) + 2ab\, C(\xi'_0, \xi_0). \tag{6.21}$$

Here, $a = \mu'_Y$ and $b = -2\mu^2_Y$, giving:

$$V(\hat{S}^2_{Y,DB}) = \eta^2\left[\mu'^2_Y V(\xi'_0) + 4\mu^4_Y V(\xi_0) - 4\mu^2_Y\mu'_Y C(\xi'_0, \xi_0)\right]. \tag{6.22}$$

Next, the covariance between $\bar{Y}'_D$ and $\bar{Y}_D$ is obtained as $C(\bar{Y}'_D, \bar{Y}_D) = \mu_Y \mu'_Y V_{11}^{(*)}(1 - \varrho\varrho')$, which accounts for the interaction between the linear and squared observations after adjustment. Finally, applying the variance decomposition formula for $\hat{S}^2_{Y,DB} = \eta(\bar{Y}'_D - (\bar{Y}_D)^2)$ and keeping only leading-order terms, we have $V_{20} = V(\xi_0)$, $V'_{20} = V(\xi'_0)$, $V_{11}^{(*)} = C(\xi'_0, \xi_0)$, $\varrho$ = correlation adjustment for $\xi_0$, and $\varrho'$ = correlation adjustment for $\xi'_0$. Then: $V(\xi_0) = V_{20}(1 - \varrho^2)$, $V(\xi'_0) = V'_{20}(1 - \varrho'^2)$, and $C(\xi'_0, \xi_0) = V_{11}^{(*)}(1 - \varrho\varrho')$.

Substituting these into the linear combination formula yields:

$$V(\hat{S}^2_{Y,DB}) = \eta^2 \left[ \mu'^2_Y V'_{20}(1 - \varrho'^2) + 4\mu^4_Y V_{20}(1 - \varrho^2) - 4\mu^2_Y \mu'_Y V_{11}^{(*)}(1 - \varrho\varrho') \right], \tag{6.23}$$

which completes the proof. $\qquad\square$

This expression shows that the variance of the difference-type estimator is reduced relative to the simple estimators due to the use of auxiliary information, while the mixed covariance term $V_{11}^{(*)}$ captures the interaction between the linear and squared observations, ensuring accurate estimation of the total variability.

**Lemma 6.7.** *Unbiased estimators of finite population variance, utilizing a single auxiliary variable X within an S-sampling scheme, are provided by:*

$$\hat{S}^2_{Y,DU} = \hat{S}^2_{Y,DB} + \eta V(\bar{Y}_D). \tag{6.24}$$

*Proof.* The demonstration can be easily achieved by following the steps in the previous proof. The variances of $\hat{S}^2_{Y,DB}$ and $\hat{S}^2_{Y,DU}$ are identical. $\qquad\square$

**Lemma 6.8.** *Variance of the $\hat{S}^2_{Y,DU}$ is given as follows:*

$$V(\hat{S}^2_{Y,DU}) \approx \eta^2 \left[ \mu'^2_Y V'_{20}(1 - \varrho'^2) + 4\mu^4_Y V_{20}(1 - \varrho^2) - 4\mu^2_Y \mu'_Y V_{11}^{(*)}(1 - \varrho\varrho') \right]. \tag{6.25}$$

*Proof.* $\hat{S}^2_{Y,DU}$ is the unbiased version of $\hat{S}^2_{Y,DB}$:

$$\hat{S}^2_{Y,DU} = \hat{S}^2_{Y,DB} + \eta V(\bar{Y}_D), \tag{6.26}$$

where $\hat{S}^2_{Y,DB} = \eta[\bar{Y}'_D - (\bar{Y}_D)^2]$ and $V(\bar{Y}_D) = \mu^2_Y V_{20}(1 - \varrho^2)$. Then the variance of a sum of two terms is:

$$V(\hat{S}^2_{Y,DU}) = V(\hat{S}^2_{Y,DB}) + V(\eta V(\bar{Y}_D)) + 2C(\hat{S}^2_{Y,DB}, \eta V(\bar{Y}_D)). \tag{6.27}$$

The term $V(\eta V(\bar{Y}_D))$ is of higher order (since $V(\bar{Y}_D)$ is already small in first-order approximation) and can be ignored. Similarly, the covariance term is negligible at the first order. Therefore, we have

$$V(\hat{S}^2_{Y,DU}) \approx V(\hat{S}^2_{Y,DB}). \tag{6.28}$$

The first-order variance still depends on the variance and covariance of linear and squared observations: $V_{20}, V'_{20}, V_{11}^{(*)}$, with the corresponding correlation reduction factors: $1 - \varrho^2$, $1 - \varrho'^2$, $1 - \varrho\varrho'$, respectively. Thus, the final variance of $\hat{S}^2_{Y,DU}$ is:

$$V(\hat{S}^2_{Y,DU}) \approx \eta^2 \left[ \mu'^2_Y V'_{20}(1 - \varrho'^2) + 4\mu^4_Y V_{20}(1 - \varrho^2) - 4\mu^2_Y \mu'_Y V_{11}^{(*)}(1 - \varrho\varrho') \right], \tag{6.29}$$

$V(\hat{S}^2_{Y,DU}) \approx V(\hat{S}^2_{Y,DB})$, which completes the proof. $\qquad\square$

Remark: The variance of the fourth estimator is identical to that of the third estimator under first-order approximation because the additional term $\eta V(\bar{Y}_D)$ is a constant (known from the sample) and contributes negligibly to the total variance at this order.

**Table 2.** Summary of notations.

| No. | Symbol / Estimator | Nomenclature | Description |
|---|---|---|---|
| 1 | $N$ | PSU | Number of clusters at stage 1 |
| 2 | $N_i$ | SSU | Number of clusters at stage 2 |
| 3 | $\overline{N}$ | Average Cluster size | Average cluster size under the 2SCS |
| 4 | $\sigma^2_{Y,b,2}$ | Between Variance | Variance between clusters under the 2SCS |
| 5 | $\sigma^2_{Y,i,2}$ | Within Variance | Variance within clusters under the 2SCS |
| 6 | $N_{ij}$ | TSU | Number of clusters at stage 3 |
| 7 | $\overline{T}$ | Average Cluster size | Average cluster size under the 3SCS |
| 8 | $\sigma^2_{Y,b,3}$ | Between Variance | Variance between clusters under the 3SCS |
| 9 | $\sigma^2_{Y,i,3}$ | Within Variance | Variance within clusters under the 3SCS |
| 10 | $S^2_Y$ | Population Variance | Variance of the study variable |
| 11 | $\hat{S}^2_{Y,SB}$ | Biased Variance Estimator | Biased estimator under S-Scheme |
| 12 | $\hat{S}^2_{Y,SU}$ | Unbiased Variance Estimator | Unbiased estimator under S-Scheme |
| 13 | $\hat{S}^2_{Y,DB}$ | Biased Difference Variance | Biased difference variance estimator |
| 14 | $\hat{S}^2_{Y,DU}$ | Unbiased Difference Variance | Unbiased difference variance estimator |
| 15 | $AB_{SB}$ | Absolute Bias SB | Absolute Bias of $\hat{S}^2_{Y,SB}$ |
| 16 | $AB_{SU}$ | Absolute Bias SU | Absolute Bias of $\hat{S}^2_{Y,SU}$ |
| 17 | $AB_{DB}$ | Absolute Bias DB | Absolute Bias of $\hat{S}^2_{Y,DB}$ |
| 18 | $AB_{DU}$ | Absolute Bias DU | Absolute Bias of $\hat{S}^2_{Y,DU}$ |
| 19 | $RE_1$ | Relative Efficiency 1 | RE of $\hat{S}^2_{Y,SB}$ with respect to $\hat{S}^2_{Y,DB}$ |
| 20 | $RE_2$ | Relative Efficiency 2 | RE of $\hat{S}^2_{Y,SB}$ with respect to $\hat{S}^2_{Y,DU}$ |
| 21 | $RE_3$ | Relative Efficiency 3 | RE of $\hat{S}^2_{Y,SU}$ with respect to $\hat{S}^2_{Y,DB}$ |
| 22 | $RE_4$ | Relative Efficiency 4 | RE of $\hat{S}^2_{Y,SU}$ with respect to $\hat{S}^2_{Y,DU}$ |

It is important to note that the terms $SU$, $SB$, $DB$, and $DU$ in the subscripts of $\hat{S}^2_Y$ represent different estimators of the finite population variance. Specifically, $SU$ and $SB$ refer to estimators that are computed without the use of auxiliary information, whereas $DB$ and $DU$ are estimated using auxiliary information. More precisely, $SU$ refers to an unbiased estimator, while $SB$ represents a biased estimator without auxiliary information. Similarly, $DB$ corresponds to a biased estimator of the finite population variance obtained through the difference estimator, whereas $DU$ denotes its unbiased counterpart derived using auxiliary information. For instance, $\hat{S}^2_{Y,SB}$ denotes the finite population variance estimator calculated without auxiliary information which is biased. Furthermore, we propose its unbiased counterpart, $\hat{S}^2_{Y,SU}$, which is also computed without the use of auxiliary information. In contrast, $\hat{S}^2_{Y,DB}$ is an estimator that utilizes auxiliary information to estimate the finite population variance which is biased. Similarly, the unbiased version of this estimator, $\hat{S}^2_{Y,DU}$, is also derived using auxiliary information. Furthermore, $AB_{SB}$, $AB_{SU}$, $AB_{DB}$, and $AB_{DU}$ denote the ABs of the corresponding variance estimators $\hat{S}^2_{Y,SB}$, $\hat{S}^2_{Y,SU}$, $\hat{S}^2_{Y,DB}$, and $\hat{S}^2_{Y,DU}$, measuring the magnitude of deviation from the true population variance. For more details, see Table 2 where a consolidated notation table is

provided for clarity.

## 6.3. Mathematical efficiency comparison between $\bar{Y}_S$, $\bar{Y}_D$, and $\bar{Y}_{Reg}$

It can be easy to show that the variance of $\bar{Y}_S$ can be written as follows [29]:

$$\text{MSE}_{\min}(\bar{Y}_S) = \mu_Y^2 \Lambda_{20},$$

and similarly, from Eq (4.5), the variance of $\bar{Y}_D$ is given as follows:

$$\text{MSE}_{\min}(\bar{Y}_D) = \mu_Y^2 \left( \Lambda_{20} - \frac{\Lambda_{11}^2}{\Lambda_{02}} \right).$$

To compare the efficiencies of the simple mean and the difference estimator, we consider the difference between their minimum MSEs. Specifically, we evaluate $\text{MSE}_{\min}(\bar{Y}) - \text{MSE}_{\min}(\bar{Y}_D)$.

$$\text{MSE}_{\min}(\bar{Y}_S) - \text{MSE}_{\min}(\bar{Y}_D) = \mu_Y^2 \Lambda_{20} - \mu_Y^2 \left( \Lambda_{20} - \frac{\Lambda_{11}^2}{\Lambda_{02}} \right)$$

$$= \mu_Y^2 \left( \Lambda_{20} - \Lambda_{20} + \frac{\Lambda_{11}^2}{\Lambda_{02}} \right) = \mu_Y^2 \frac{\Lambda_{11}^2}{\Lambda_{02}}.$$

Since $\mu_Y^2 > 0$ and (under the usual regularity conditions) $\Lambda_{02} > 0$, the difference is nonnegative and equals zero if and only if $\Lambda_{11} = 0$. Thus, the minimum MSE of the difference estimator is smaller (or equal, in the degenerate case) than that of the simple mean by the amount $\mu_Y^2 \Lambda_{11}^2 / \Lambda_{02}$. In the case of the regression estimator, for large samples, $\hat{\beta}_S$ is a consistent estimator of the true $\beta_S$, since as the sample size increases, $\hat{\beta}_S$ converges to its true value. On this basis, the performance of the regression estimator improves in large samples, and its bias becomes negligible. Consequently, the minimum MSE of the regression estimator as given in Eq 4.6 ultimately equals that of the difference estimator, implying that, asymptotically, both estimators perform equivalently.

The schematic diagram in Figure 1 illustrates the hierarchical link between the 2SCS and 3SCS. In the 2SCS framework, the population is first divided into PSUs, and subsequently, SSUs are selected within the chosen PSUs. In contrast, the 3SCS framework introduces an additional level of sampling, whereby TSUs are further selected within each chosen SSU. From a methodological perspective, the 3SCS formulation generalizes that of the 2SCS: when the third stage is omitted (i.e., no TSUs are drawn), the MSE and bias expressions for the 3SCS reduce directly to those of the 2SCS. Thus, the two-stage design can be regarded as a special case of the more general 3SCS framework.
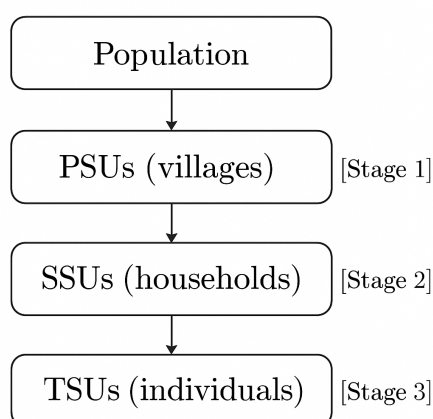
**Figure 1.** Schematic link showing the 2SCS as a special case of the generalized 3SCS framework.

In the next section, we validate the proposed estimators using two real-life datasets. Specifically, the biased and unbiased estimators constructed with auxiliary information are compared without auxiliary information. The comparison is based on the MSE, which is subsequently used to compute the RE of the estimators. This choice is justified, as comparing MSEs provides an appropriate measure of both bias and variability, allowing for a meaningful assessment of estimator performance.

## 7. Empirical study

This section analyzes real datasets to determine the AB, MSE, and RE of the proposed variance estimators. These estimators are derived from the 2SCS and 3SCS methods, both with and without the inclusion of auxiliary information.

### 7.1. Population I

Population I is drawn from a cross-sectional survey conducted in the Multan District, Pakistan (Jan–Mar 2020) [36]. The sampling frame consisted of school-going children aged 3–18 years, enrolled in government and private schools, with a total population of 4.7 million residents in the district. Schools were first stratified by socioeconomic status (SES) based on fee structure, forming the PSUs. Within selected schools, classes served as the SSUs, and for 3-stage sampling, students were further grouped into age categories (TSUs). In total, 1040 students were selected using a probability proportional to size (PPS). The body mass index (BMI), calculated as weight (kg)/height (m$^2$), was used as the study variable, while body weight (kg) served as the auxiliary variable due to its strong correlation with the BMI. In the original study, written informed consent was obtained from the school administration, and verbal consent was taken from students prior to participation. For empirical analysis, the dataset was reorganized into two strata (males and females), while SES remained the PSU. Thus, for the 2SCS design, SES and BMI were treated as the PSU and SSU, respectively, whereas for the 3SCS design, SES, age, and BMI were considered as the PSU, SSU, and TSU, respectively.

## 7.2. Population II

The other dataset is sourced from the Centers for Disease Control and Pervention is associated with the second national health and nutrition examination survey. The dataset comprises 10,351 units representing the non-institutionalized civilian population of the United States, including all 50 states and the District of Columbia. The data were divided into four geographic regions (REGs)—midwestern, southern, northeastern, and western—each further subdivided into specific locations (LOCs), with REG1, REG2, REG3, and REG4 containing 16, 14, 16, and 16 LOCs, respectively. The total number of units in all LOCs of REG1 is 2774 (average LOC size: 173.38); in REG2, 2096 (average: 149.71); in REG3, 2853 (average: 178.31); and in REG4, 2628 (average: 164.25). Random stratification into two strata was performed using a Bernoulli distribution with success probability 0.50, where 0 indicates Stratum I and 1 indicates Stratum II. The study variable $Y$ represents BMI, and the auxiliary variable $X$ is body weight (kg), which is highly correlated with BMI. For the 2SCS design, LOC and BMI are treated as the PSU and SSU, respectively, whereas for the 3SCS design, REG, LOC, and BMI are considered the PSU, SSU, and TSU, respectively. The dataset is publicly available and can be downloaded from `https://www.stata-press.com/data/r15/svy.html` or via Figshare at `https://doi.org/10.6084/m9.figshare.28151576.v1`, where it has been uploaded for public availability.

## 7.3. Numerical illustration

Using two real datasets, we evaluate the performance of the proposed variance estimators under the 2SCS and 3SCS schemes. For this purpose, we compute the AB, MSE, and RE of the estimators under an S-sampling scheme by using the following formulas:

$$
\begin{aligned}
\mathrm{AB}(\hat{S}^2_{Y,SB}) &= \left| \frac{1}{r} \sum_{i=1}^{r} \hat{S}^2_{Y,SB,i} - S^2_Y \right|, \\
\mathrm{MSE}(\hat{S}^2_{Y,SB}) &= \frac{1}{r} \sum_{i=1}^{r} \left( \hat{S}^2_{Y,SB,i} - S^2_Y \right)^2, \\
\mathrm{MSE}(\hat{S}^2_{Y,DB}) &= \frac{1}{r} \sum_{i=1}^{r} \left( \hat{S}^2_{Y,DB,i} - S^2_Y \right)^2.
\end{aligned}
$$

To obtain stable estimates of these measures, we iterate the computations 10,000 times (i.e., $r = 10,000$). The ABs and MSEs of these estimators were estimated for different $n$, $n_i$, and $n_{ij}$ as provided in Tables 3–6. The ABs and MSEs of the biased estimators with and without auxiliary information are calculated using the formulas provided above. The RE of the biased estimator of $S^2_Y$ with respect to the unbiased variance estimator is given by

$$
\mathrm{RE}(\hat{S}^2_{Y,SB}, \hat{S}^2_{Y,DB}) = \frac{\mathrm{MSE}(\hat{S}^2_{Y,SB})}{\mathrm{MSE}(\hat{S}^2_{Y,DB})}. \tag{7.1}
$$

**Table 3.** ABs and REs of the proposed estimators with respect to $\hat{S}^2_{Y,SB}$ under the 2SCS using Population I.

| $n$ | $n_i$ | $AB_{SB}$ | $AB_{SU}$ | $AB_{DB}$ | $AB_{DU}$ | $RE_1$ | $RE_2$ | $RE_3$ | $RE_4$ |
|---|---|---|---|---|---|---|---|---|---|
| 2 | 03 | 14.09 | 1.66 | 2.03 | 0.26 | 07.72 | 07.79 | 07.61 | 07.67 |
| 2 | 05 | 12.27 | 0.81 | 2.06 | 0.55 | 10.75 | 10.88 | 10.76 | 10.89 |
| 2 | 10 | 13.32 | 2.28 | 0.77 | 0.05 | 22.78 | 22.88 | 22.61 | 22.70 |
| 2 | 15 | 09.58 | 0.77 | 0.19 | 0.38 | 45.65 | 45.58 | 45.94 | 45.87 |
| 2 | 20 | 10.12 | 0.37 | 0.09 | 0.33 | 62.61 | 62.49 | 63.20 | 63.08 |
| 3 | 03 | 04.35 | 0.32 | 1.05 | 0.21 | 04.00 | 04.02 | 03.98 | 03.99 |
| 3 | 05 | 02.61 | 1.71 | 0.06 | 0.81 | 05.83 | 05.80 | 05.82 | 05.80 |
| 3 | 10 | 02.27 | 1.46 | 0.50 | 0.00 | 12.28 | 12.31 | 12.30 | 12.34 |
| 3 | 15 | 02.08 | 1.54 | 0.26 | 0.09 | 17.66 | 17.68 | 17.68 | 17.70 |
| 3 | 20 | 03.56 | 0.01 | 0.25 | 0.02 | 24.95 | 24.99 | 24.89 | 24.93 |

Note: $AB_{SB}$, $AB_{SU}$, $AB_{DB}$, and $AB_{DU}$ denote the AB of $\hat{S}^2_{Y,SB}$ $\hat{S}^2_{Y,SU}$, $\hat{S}^2_{Y,DB}$, and $\hat{S}^2_{Y,DU}$, respectively.

**Table 4.** ABs and REs of the proposed estimators with respect to $\hat{S}^2_{Y,SB}$ under the 2SCS using Population II.

| $n$ | $ni$ | $AB_{SB}$ | $AB_{SU}$ | $AB_{DB}$ | $AB_{DU}$ | $RE_1$ | $RE_2$ | $RE_3$ | $RE_4$ |
|---|---|---|---|---|---|---|---|---|---|
| 3 | 05 | 18.23 | 2.54 | 0.48 | 0.18 | 085.19 | 085.33 | 078.84 | 078.96 |
| 3 | 10 | 14.68 | 0.01 | 0.34 | 0.01 | 167.49 | 167.82 | 156.25 | 156.56 |
| 3 | 15 | 11.29 | 3.71 | 0.07 | 0.31 | 221.87 | 221.39 | 208.28 | 207.83 |
| 3 | 20 | 15.06 | 0.05 | 0.21 | 0.03 | 386.50 | 387.14 | 364.78 | 365.39 |
| 5 | 05 | 08.79 | 0.70 | 0.43 | 0.04 | 073.24 | 073.42 | 069.67 | 069.84 |
| 5 | 10 | 10.66 | 1.67 | 0.18 | 0.02 | 166.55 | 166.72 | 155.59 | 155.74 |
| 5 | 15 | 09.68 | 0.84 | 0.11 | 0.03 | 232.55 | 232.66 | 221.41 | 221.51 |
| 5 | 20 | 10.21 | 1.43 | 0.12 | 0.23 | 317.67 | 317.00 | 302.58 | 301.94 |
| 10 | 05 | 03.90 | 0.47 | 0.34 | 0.15 | 070.72 | 070.92 | 069.33 | 069.53 |
| 10 | 10 | 05.77 | 1.65 | 0.25 | 0.14 | 122.01 | 122.27 | 118.39 | 118.64 |
| 10 | 15 | 02.79 | 1.19 | 0.18 | 0.11 | 207.60 | 208.00 | 203.12 | 203.50 |
| 10 | 20 | 04.19 | 0.21 | 0.08 | 0.03 | 272.96 | 273.16 | 265.92 | 266.11 |

**Table 5.** ABs and REs of the proposed estimators with respect to $\hat{S}^2_{Y,SB}$ under the 3SCS using Population I.

| $n$ | $n_i$ | $n_{ij}$ | $AB_{SB}$ | $AB_{SU}$ | $AB_{DB}$ | $AB_{DU}$ | $RE_1$ | $RE_2$ | $RE_3$ | $RE_4$ |
|---|---|---|---|---|---|---|---|---|---|---|
| 2 | 3 | 3 | 30.46 | 7.75 | 1.55 | 0.32 | 100.43 | 101.97 | 083.75 | 085.03 |
| 2 | 3 | 4 | 30.00 | 7.92 | 1.33 | 0.22 | 108.53 | 109.90 | 091.20 | 092.36 |
| 2 | 3 | 5 | 29.72 | 7.56 | 1.26 | 0.22 | 109.23 | 110.48 | 091.92 | 092.97 |
| 2 | 5 | 3 | 24.57 | 8.71 | 0.95 | 0.21 | 106.75 | 107.67 | 093.27 | 094.07 |
| 2 | 5 | 4 | 25.07 | 9.43 | 0.88 | 0.22 | 123.95 | 124.99 | 108.76 | 109.67 |
| 2 | 5 | 5 | 23.22 | 7.65 | 0.92 | 0.31 | 131.52 | 132.72 | 115.98 | 117.04 |
| 2 | 7 | 3 | 20.23 | 7.37 | 0.75 | 0.25 | 124.67 | 125.55 | 112.78 | 113.57 |
| 2 | 7 | 4 | 21.21 | 8.22 | 0.70 | 0.25 | 138.68 | 139.61 | 124.58 | 125.42 |
| 2 | 7 | 5 | 20.41 | 7.85 | 0.60 | 0.19 | 153.03 | 153.89 | 138.24 | 139.02 |
| 3 | 3 | 3 | 20.40 | 7.98 | 0.86 | 0.12 | 085.68 | 086.31 | 075.34 | 075.90 |
| 3 | 3 | 4 | 20.80 | 8.75 | 0.82 | 0.15 | 100.19 | 100.94 | 088.76 | 089.42 |
| 3 | 3 | 5 | 19.07 | 7.15 | 0.91 | 0.29 | 098.14 | 098.98 | 087.09 | 087.84 |
| 3 | 5 | 3 | 16.56 | 8.65 | 0.58 | 0.14 | 097.34 | 097.81 | 088.99 | 089.42 |
| 3 | 5 | 4 | 15.51 | 7.83 | 0.69 | 0.29 | 111.10 | 111.84 | 101.96 | 102.64 |
| 3 | 5 | 5 | 16.21 | 8.54 | 0.54 | 0.17 | 125.78 | 126.42 | 115.14 | 115.73 |
| 3 | 7 | 3 | 13.07 | 7.14 | 0.53 | 0.23 | 108.86 | 109.38 | 101.79 | 102.27 |
| 3 | 7 | 4 | 13.71 | 7.91 | 0.40 | 0.13 | 125.65 | 126.07 | 117.62 | 118.01 |
| 3 | 7 | 5 | 13.10 | 7.42 | 0.31 | 0.07 | 140.70 | 141.04 | 132.37 | 132.68 |

**Table 6.** ABs and REs of the proposed estimators with respect to $\hat{S}^2_{Y,SB}$ under the 3SCS using Population II.

| $n$ | $n_i$ | $n_{ij}$ | $AB_{SB}$ | $AB_{SU}$ | $AB_{DB}$ | $AB_{DU}$ | $RE_1$ | $RE_2$ | $RE_3$ | $RE_4$ |
|---|---|---|---|---|---|---|---|---|---|---|
| 2 | 3 | 05 | 12.60 | 2.37 | 0.47 | 0.14 | 118.42 | 118.81 | 114.38 | 114.76 |
| 2 | 3 | 10 | 12.25 | 2.84 | 0.24 | 0.07 | 241.10 | 241.52 | 234.49 | 234.90 |
| 2 | 3 | 15 | 12.24 | 3.17 | 0.06 | 0.06 | 368.88 | 368.89 | 357.46 | 357.46 |
| 2 | 5 | 05 | 8.24 | 1.29 | 0.30 | 0.10 | 134.25 | 134.53 | 133.92 | 134.20 |
| 2 | 5 | 10 | 9.45 | 3.16 | 0.07 | 0.03 | 292.28 | 292.34 | 290.77 | 290.83 |
| 2 | 5 | 15 | 8.95 | 2.92 | 0.09 | 0.02 | 443.40 | 443.68 | 442.20 | 442.47 |
| 2 | 7 | 05 | 7.71 | 2.19 | 0.22 | 0.08 | 165.38 | 165.65 | 166.66 | 166.94 |
| 2 | 7 | 10 | 7.63 | 2.60 | 0.08 | 0.01 | 335.87 | 336.02 | 339.36 | 339.52 |
| 2 | 7 | 15 | 6.43 | 1.56 | 0.09 | 0.04 | 520.34 | 520.76 | 527.26 | 527.69 |
| 3 | 3 | 05 | 8.40 | 2.29 | 0.25 | 0.03 | 111.17 | 111.33 | 107.42 | 107.57 |
| 3 | 3 | 10 | 8.76 | 3.32 | 0.07 | 0.04 | 234.78 | 234.82 | 227.32 | 227.36 |
| 3 | 3 | 15 | 7.14 | 1.95 | 0.09 | 0.01 | 340.86 | 341.05 | 331.71 | 331.89 |
| 3 | 5 | 05 | 6.27 | 2.33 | 0.18 | 0.05 | 128.63 | 128.77 | 126.39 | 126.53 |
| 3 | 5 | 10 | 5.57 | 2.15 | 0.09 | 0.02 | 269.82 | 269.97 | 266.57 | 266.72 |
| 3 | 5 | 15 | 7.01 | 3.79 | 0.05 | 0.01 | 418.16 | 418.30 | 412.51 | 412.65 |
| 3 | 7 | 05 | 4.65 | 1.65 | 0.13 | 0.04 | 150.00 | 150.13 | 148.91 | 149.03 |
| 3 | 7 | 10 | 5.18 | 2.65 | 0.07 | 0.02 | 313.80 | 313.95 | 312.04 | 312.20 |
| 3 | 7 | 15 | 5.04 | 2.67 | 0.04 | 0.00 | 483.15 | 483.25 | 481.22 | 481.33 |

On similar lines, one can find the ABs, MSEs, and REs of the other aforementioned estimators, as given in Section 6. It should be noted that from Tables 3–6, $RE_1$, $RE_2$, $RE_3$, and $RE_4$ correspond to the RE of $\hat{S}^2_{Y,SB}$ with $\hat{S}^2_{Y,DB}$, $\hat{S}^2_{Y,SB}$ with $\hat{S}^2_{Y,DU}$, $\hat{S}^2_{Y,SU}$ with $\hat{S}^2_{Y,DB}$, and $\hat{S}^2_{Y,SU}$ with $\hat{S}^2_{Y,DU}$, respectively. All numerical calculations and iteration have been carried out using R 4.4.1.

To validate our proposed variance estimators, the true variance of Population II ($S^2_Y = 10.8022$) was computed using Eq (5.2). Under the 3SCS scheme, with $n = 3$, $n_i = 7$, and $n_{ij} = 5$, the estimated variances obtained are $\hat{S}^2_{Y,SB} = 2.28932$, $\hat{S}^2_{Y,SU} = 3.427977$, $\hat{S}^2_{Y,DB} = 10.40159$, and $\hat{S}^2_{Y,DU} = 10.64815$. This validation was conducted only for the 3SCS, while the 2SCS scheme was omitted to avoid repetition, as similar conclusions are expected in both designs. It can be observed that all the proposed estimators provide estimates of the true variance; however, the difference estimators ($\hat{S}^2_{Y,DB}$ and $\hat{S}^2_{Y,DU}$) are much closer to the true variance, thereby confirming their superior performance in practical applications.

For $n = 2$, $AB_{SB}$ and $AB_{SU}$ show moderate variation across increasing $n_i$, while $AB_{DB}$ and $AB_{DU}$ decrease notably, with $AB_{DU}$ nearing zero at $n_i = 10$ as we see in Figure 2. The REs ($RE_1$ to $RE_4$) increase steadily, peaking around 63 at $n = 2$, $n_i = 20$, highlighting improved estimator precision.
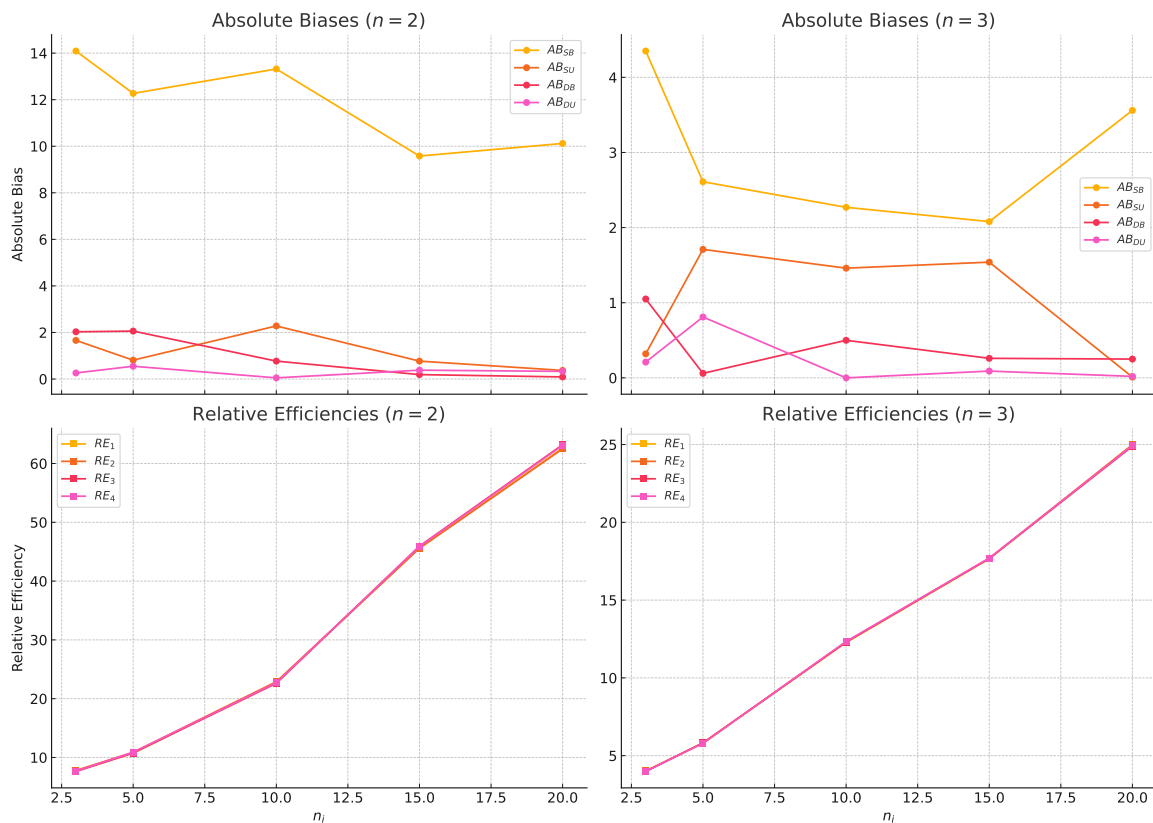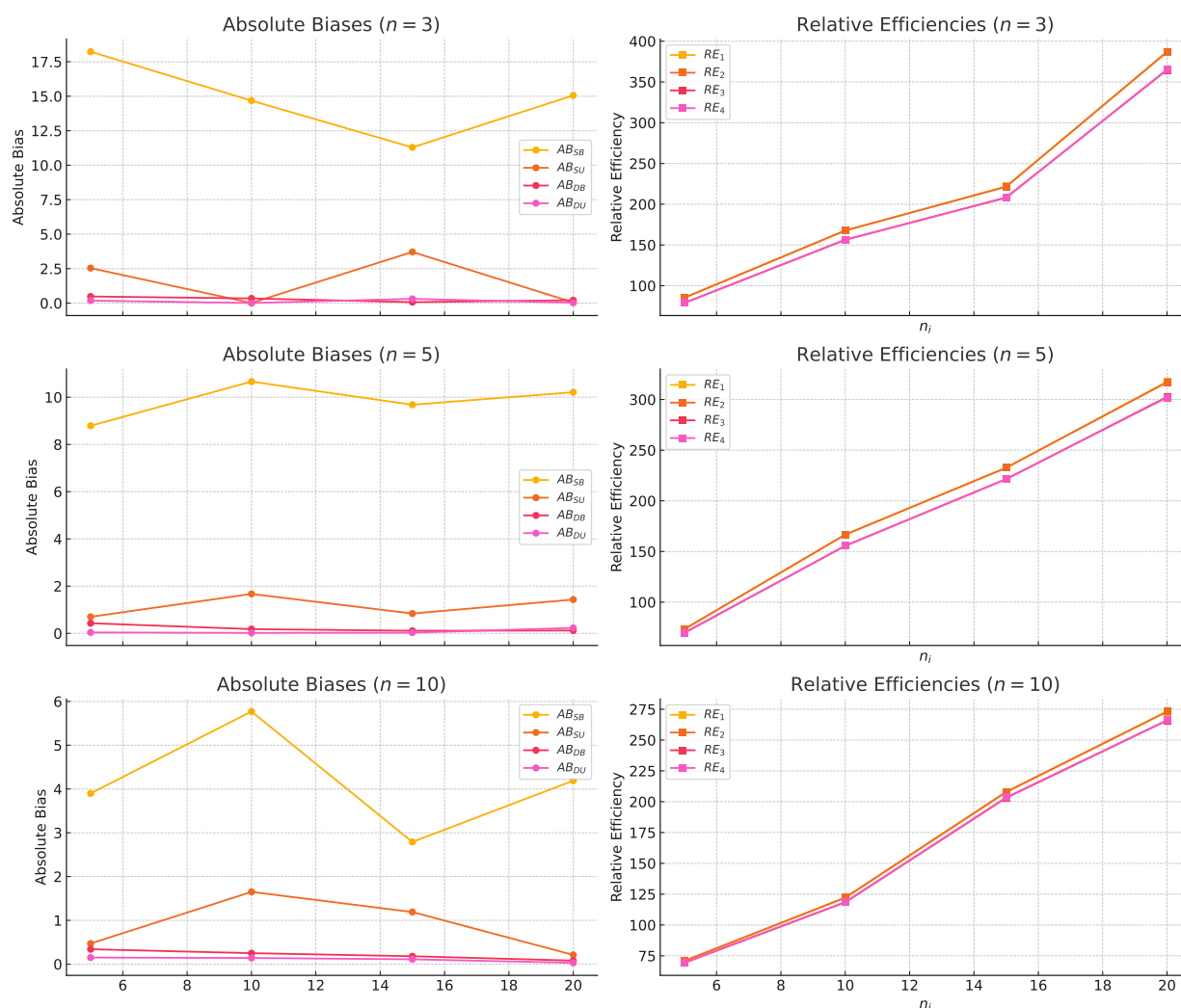
**Figure 2.** ABs and REs of the proposed estimators with respect to $\hat{S}^2_{Y,SB}$ under the 2SCS using Population I.

When $n = 3$, all ABs are lower, and REs remain consistently high, with DB and DU estimators performing best in terms of accuracy and efficiency.

In Population II (Table 4), the estimators demonstrate greater improvement with increased $n$ and $n_i$. $AB_{\mathrm{DU}}$ and $AB_{\mathrm{DB}}$ remain minimal across all cases, while REs rise sharply, exceeding 380 for $n = 3$, $n_i = 20$, and which can also be seen in Figure 3. This confirms that deeper cluster sampling significantly enhances precision, especially for DU and DB estimators.

**Figure 3.** ABs and REs of the proposed estimators with respect to $\hat{S}^2_{Y,SB}$ under the 2SCS using Population II.

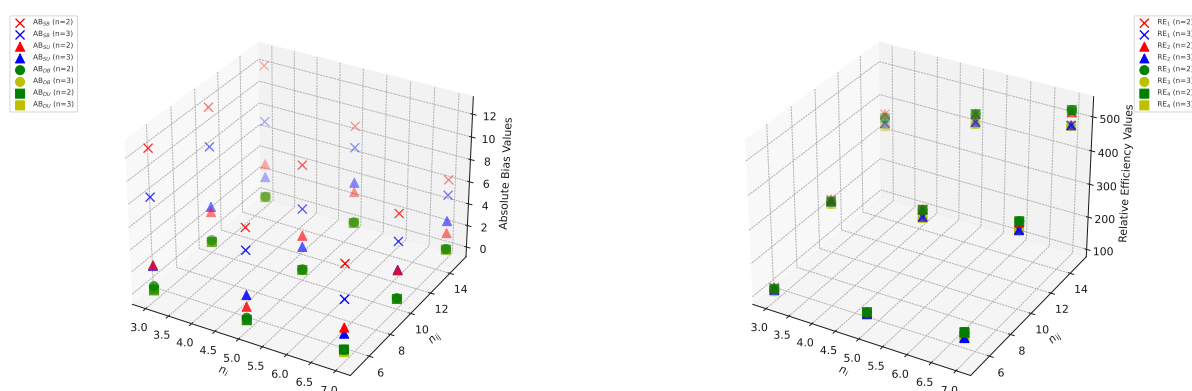The 3D scatter plots given in Figures 4 and 5 using the values given in Tables 5 and 6, respectively, illustrate how ABs ($AB_{SB}, AB_{SU}, AB_{DB}, AB_{DU}$) and REs ($RE_1$ to $RE_4$) of the proposed estimators vary with $n_i$ and $n_{ij}$ for $n = 2$ and $n = 3$. As both $n_i$ and $n_{ij}$ increase, the ABs generally decrease, indicating improved estimator accuracy, while the REs increase, reflecting better performance with larger sub-sample sizes.

**Figure 4.** ABs and REs of the proposed estimators with respect to $\hat{S}^2_{Y,SB}$ under the 3SCS using Population I.



**Figure 5.** ABs and REs of the proposed estimators with respect to $\hat{S}^2_{Y,SB}$ under the 3SCS using Population II

Notably, $AB_{\text{DU}}$ and $AB_{\text{DB}}$ exhibit the lowest bias, while $RE_3$ and $RE_4$ consistently achieve higher efficiency, especially for $n = 3$. Under the 3SCS using Population I, the difference unbiased estimator $\hat{S}^2_{Y,DU}$ attained the highest efficiency ($RE_2 = 153.89$), closely followed by the difference biased estimator $\hat{S}^2_{Y,DB}$ ($RE_1 = 527.26$). Similarly, using Population II, $\hat{S}^2_{Y,DU}$ attained the highest efficiency ($RE_4 = 527.69$), closely followed by the $\hat{S}^2_{Y,DB}$ ($RE_3 = 527.26$). Both substantially outperformed the conventional sample variance estimator without auxiliary information. These trends suggest that increasing the complexity of the configuration, through larger $n_i$ and $n_{ij}$ values, leads to lower bias and higher efficiency, enhancing overall estimator performance.

### 7.4. Robustness and limitations for small samples

We examined the performance of the proposed estimators under very small sample sizes (e.g., $n = 1, n_i = 2, n_{ij} = 5$) to assess robustness. Some variance estimates could not be computed (N/A) because multi-stage cluster sampling requires a minimum number of clusters and subunits for valid estimation. Nonetheless, under the 3SCS scheme using Population II, the difference estimators yielded reasonable estimates: $\hat{S}^2_{Y,DB} = 4.6539$, $\hat{S}^2_{Y,DU} = 10.8093$, with ABs $AB_{DB} = 6.14834$, $AB_{DU} = 0.0071$, and REs $RE_1 = 177.36$, $RE_2 = 186.72$. These results indicate that while the difference estimators are relatively robust even in sparse data settings, the accuracy and stability of variance estimation deteriorate when the number of clusters or subunits is extremely small. Therefore, we recommend a minimum number of clusters and subunits to ensure reliable estimation in practical applications.

The RE of the proposed estimators increases significantly with larger sample sizes, as higher $n$, $n_i$, and $n_{ij}$ reduce AB and provide more precise variance estimates. For survey practitioners, this indicates that increasing sub-sample sizes and using auxiliary information can substantially improve estimator performance. However, auxiliary information may fail to improve efficiency if it is poorly correlated with the study variable, measured with error, or has incorrect distributional assumptions, potentially increasing bias or variance. Moreover, estimators using auxiliary information can be sensitive to outliers or skewed auxiliary variables, which may reduce efficiency or introduce bias. Among the estimators considered, the $\hat{S}^2_{Y,DU}$ consistently performs best, exhibiting the lowest bias and highest efficiency across 2SCS and 3SCS schemes. While regression estimators could be applied, their efficiency is generally lower due to the need to estimate regression coefficients from the sample. These findings highlight practical strategies for designing surveys that yield reliable and robust variance estimates. Additionally, it is important to consider uncertainties arising from non-sampling errors, measurement inaccuracies, and assumptions regarding auxiliary information, as these factors may impact the reliability of the estimators in real-world applications.

### 7.5. Discussion

In this study, we have estimated the finite population variance under the 2SCS and 3SCS using the method of moments. Accurate and efficient variance estimation enhances the reliability of survey-based inferences, leading to more precise estimates, better resource allocation, increased confidence in policy decisions, effective monitoring and evaluation of programs, and identification of priority populations. These advantages support informed decision-making in health and demographic policy contexts. For this purpose, variance estimation has been carried out through both the difference estimator and the simple estimator (without auxiliary information), and their results have been compared.

The numerical results reported in Tables 3–6 consistently demonstrate that the difference unbiased estimator ($\hat{S}^2_{Y,DU}$) outperforms the other proposed and conventional estimators under both the 2SCS and 3SCS schemes. This superior performance is fully consistent with the theoretical expectations, since the estimator leverages the structure of the difference estimator together with raw moment information, leading to both lower AB and higher RE. The maximum efficiency observed under the 3SCS ($RE_3 = 527.69$ for $\hat{S}^2_{Y,DU}$) validates the theoretical prediction that difference-type estimators incorporating auxiliary information should outperform sample-based estimators without auxiliary variables. Thus, the empirical evidence not only supports but also strengthens the theoretical justification of $\hat{S}^2_{Y,DU}$ as the most reliable estimator of finite population variance.

While the proposed estimators show clear efficiency gains, certain limitations should be acknowledged. First, the numerical evaluation is based on two real-life datasets, which may constrain the generalizability of the findings to populations with different structures or auxiliary relationships. Second, although we have now derived the closed-form variance expressions for the proposed estimators under the 2SCS and 3SCS schemes (which were not provided in the work of [24] under stratified random sampling), the efficiency assessment in this study has been limited to AB, MSE, and RE. Third, the efficiency of estimators involving auxiliary information depends critically on the quality and strength of the correlation between the study and auxiliary variables; in practice, poor auxiliary data may limit their effectiveness. Therefore, while the results strongly favor the use of difference-type estimators, particularly the unbiased version, care must be taken when applying them in real-world surveys where auxiliary information may be weak or unreliable.

Implementing the proposed estimators in real surveys may face challenges such as limited auxiliary information, non-response, and computational complexity for multi-stage and stratified designs. Their efficiency depends on assumed relationships between study and auxiliary variables, which may not always hold, and uncertainties from non-sampling errors, measurement inaccuracies, and model assumptions can further affect reliability. Although our simulations using two real-life datasets account for sampling variability, they do not capture all practical uncertainties, and full uncertainty modeling is beyond the current scope, representing a direction for future research. Frameworks such as those in [37–41] could be adapted in future studies to enhance the robustness of variance estimation under complex sampling. Nevertheless, the proposed estimators rely on simple computations involving sample means, variances, covariances, and available auxiliary information, ensuring computational efficiency for typical survey applications. For large-scale surveys, the primary computational effort arises from multi-stage clustering and processing auxiliary data, which can be handled efficiently using standard statistical software, making the estimators practically feasible and providing a flexible framework for variance estimation in complex survey designs.

## 8. Conclusions

This study proposed biased and unbiased estimators of finite population variance under the 2SCS and 3SCS schemes, including difference-type estimators that exploit auxiliary information. Analytical and empirical results showed that incorporating auxiliary variables reduces bias and MSE while substantially improving RE. Notably, the proposed unbiased and biased difference estimators under the 2SCS/3SCS outperformed proposed conventional estimators (without auxiliary information), confirming their practical advantage. However, the analysis was limited to two real populations and a single auxiliary variable. Future research may extend these methods to multiple auxiliaries and other complex designs such as stratified multi-stage sampling. For survey practitioners, the findings underscore that even biased estimators leveraging auxiliary data can deliver significant efficiency gains, enabling more precise variance estimation without expanding the sample size.

### 8.1. Future works

The accuracy of the proposed estimators for finite population variance can be enhanced by incorporating multiple auxiliary variables. In future research, the difference estimator can be extended by incorporating two or more auxiliary variables (dual difference estimator, ratio-cum-

product or exponential ratio-cum-product type estimators), enabling the estimator to exploit additional information for more accurate and less biased variance estimation. Moreover, difference estimators or other classes of estimators provided by [42, 43] based on single or dual auxiliary information can also be explored to improve efficiency further. Additionally, adopting PPS sampling at the first stage would allow larger units a higher probability of selection, thereby making the sample more representative. The integration of these approaches is expected to not only enhance the precision of estimation but also provide more effective control of variance in complex survey designs. In a similar manner, new estimators for finite population variance can be developed by incorporating auxiliary information into the frameworks of stratified 2SCS and 3SCS schemes. If the estimators for the 2SCS and 3SCS are applied in stratified designs, the population must be divided into strata, and separate estimates and variances need to be calculated for each stratum, which are then combined using stratum weights [29, 33]. The bias and MSE formulas for ratio-, product-, or moment-based estimators must also be adjusted stratum-wise, and sample allocation can be optimized based on stratum size and variability [24, 44].

## Author contributions

Mohsin Abbas: Writing-original draft, writing-review and editing, conceptualization, methodology, and formal analysis. Muhammad Ahmed Shehzad: Supervision, investigation, formal analysis, and writing-review and editing. Hasnain Iftikhar: Writing-review and editing, investigation, formal analysis, data curation, resources, and project administration. Paulo Canas Rodrigues: Writing-review and editing, investigation, formal analysis, resources, and funding acquisition. Abdulmajeed Atiah Alharbi and Jeza Allohibi: Writing-review and editing, software, investigation,, visualization, formal analysis, data curation, and resources.

## Use of AI tools declaration

The authors declare that they have not used Artificial Intelligence (AI) tools in creating this article.

## Conflict of interest

The authors declare no conflicts of interest.

## References

1. M. H. Hansen, W. N. Hurwitz, On the theory of sampling from finite populations, *Ann. Math. Statist.*, **14** (1943), 333–362. http://dx.doi.org/10.1214/aoms/1177731356

2. W. G. Cochran, *Sampling techniques, 3rd Edition*, John Wiley, 1977.

3. C. E. Särndal, B. Swensson, J. Wretman, *Model assisted survey sampling*, Springer Science & Business Media, 2003.

4. J. Rao, J. Kovar, H. Mantel, On estimating distribution functions and quantiles from survey data using auxiliary information, *Biometrika*, **77** (1990), 365–375. https://doi.org/10.2307/2336815

5. R. K. Burdick, R. L. Sielken Jr, Variance estimation based on a superpopulation model in two-stage sampling, *J. Am. Stat. Assoc.*, **74** (1979), 438–440. https://doi.org/10.1080/01621459.1979.10482533

6. R. M. Royall, W. G. Cumberland, Variance estimation in finite population sampling, *J. Am. Stat. Assoc.*, **73** (1978), 351–358. https://doi.org/10.1080/01621459.1978.10481581

7. D. Haziza, J. Rao, Variance estimation in two-stage cluster sampling under imputation for missing data, *J. Stat. Theory Pract.*, **4** (2010), 827–844. https://doi.org/10.1080/15598608.2010.10412021

8. P. V. Sukhatme, B. Sukhatme, S. Sukhatme, C. Asok, *Sampling theory of surveys with applications*, 1984.

9. L. Sahoo, A regression-type estimator in two-stage sampling, *Calcutta Stat. Assoc. Bull.*, **36** (1987), 97–100. https://doi.org/10.1177/000806831987011

10. R. Yang, H. Li, H. Huang, Multisource information fusion considering the weight of focal element's beliefs: a gaussian kernel similarity approach, *Meas. Sci. Technol.*, **35** (2024), 025136. https://doi.org/10.1088/1361-6501/ad0e3b

11. S. Ahmad, H. Iftikhar, M. Qureshi, I. Khan, A. S. Omer, E. A. T. Armas, et al., A new auxiliary variables-based estimator for population distribution function under stratified random sampling and non-response, *Sci. Rep.*, **15** (2025), 13580. https://doi.org/10.1038/s41598-025-98246-y

12. M. R. Garcia, A. A. Cebrian, Repeated substitution method: The ratio estimator for the population variance, *Metrika*, **43** (1996), 101–105. https://doi.org/10.1007/BF02613900

13. L. N. Upadhyaya, H. P. Singh, S. Singh, A class of estimators for estimating the variance of the ratio estimator, *J. Japan Stat. Soc.*, **34** (2004), 47–63. https://doi.org/10.14490/jjss.34.47

14. P. Chandra, H. Singh, A family of estimators for population variance using knowledge of kurtosis of an auxiliary variable in sample survey, *Stat. Transit.*, **7** (2005), 27–34.

15. A. Arcos, M. Rueda, M. Martınez, S. González, Y. Roman, Incorporating the auxiliary information available in variance estimation, *Appl. Math. Comput.*, **160** (2005), 387–399. https://doi.org/10.1016/j.amc.2003.11.010

16. C. Kadilar, H. Cingi, Improvement in variance estimation in simple random sampling, *Commun. Stat. Theor. M.*, **36** (2007), 2075–2081. https://doi.org/10.1080/03610920601144046

17. L. K. Grover, A correction note on improvement in variance estimation using auxiliary information, *Commun. Stat. Theor. M.*, **39** (2010), 753–764. https://doi.org/10.1080/03610920902785786

18. Z. Zhou, Y. Wang, X. Liu, Z. Li, M. Wu, G. Zhou, Hybrid of neural network and physics-based estimator for vehicle longitudinal dynamics modeling using limited driving data, *IEEE T. Intell. Transp.*, 2025. https://doi.org/10.1109/TITS.2025.3585346

19. P. Sharma, R. Singh, A generalized class of estimators for finite population variance in presence of measurement errors, *J. Mod. Appl. Stat. Meth.*, **12** (2013), 231–241. https://doi.org/10.22237/jmasm/1383279120

20. H. P. Singh, R. S. Solanki, Improved estimation of finite population variance using auxiliary information, *Commun. Stat. Theor. M.*, **42** (2013), 2718–2730. https://doi.org/10.1080/03610926.2011.617485

21. S. Ahmad, N. K. Adichwal, M. Aamir, J. Shabbir, N. Alsadat, M. Elgarhy, et al., An enhanced estimator of finite population variance using two auxiliary variables under simple random sampling, *Sci. Rep.*, **13** (2023), 21444. https://doi.org/10.1038/s41598-023-44169-5

22. S. Bhushan, A. Kumar, A. Alsubie, S. A. Lone, Variance estimation under an efficient class of estimators in simple random sampling, *Ain Shams Eng. J.*, **14** (2023), 102012. https://doi.org/10.1016/j.asej.2022.102012

23. U. Daraz, M. A. Alomair, O. Albalawi, A. S. Al Naim, New techniques for estimating finite population variance using ranks of auxiliary variable in two-stage sampling, *Mathematics*, **12** (2024), 2741. https://doi.org/10.3390/math12172741

24. A. Haq, M. Usman, M. Khan, Estimation of finite population variance under stratified random sampling, *Commun. Stat. Simul. C.*, **52** (2023), 6193–6209. https://doi.org/10.1080/03610918.2021.2009866

25. A. Kumar, R. Suhail, S. Katara, Enhanced estimation of population variance under simple random sampling with an application to real data, *Int. J. Agric. Stat. Sci.*, **20** (2024), 87–96. https://doi.org/10.59467/IJASS.2024.20.87

26. S. Ahmad, M. Qureshi, H. Iftikhar, P. C. Rodrigues, M. Z. Rehman, An improved family of unbiased ratio estimators for a population distribution function, *AIMS Math.*, **10** (2025), 1061–1084. https://doi.org/10.3934/math.2025051

27. H. A. Lone, R. Tailor, Estimation of population variance in simple random sampling, *J. Stat. Manag. Syst.*, **20** (2017), 17–38. https://doi.org/10.1080/09720510.2016.1187923

28. A. Sanaullah, I. Niaz, J. Shabbir, I. Ehsan, A class of hybrid type estimators for variance of a finite population in simple random sampling, *Commun. Stat. Simul. C.*, **51** (2022), 5609–5619. https://doi.org/10.1080/03610918.2020.1776873

29. M. Abbas, M. Ahmed Shehzad, H. Khurram, M. Rabia, Estimation of finite population mean in a complex survey sampling, *PLOS One*, **20** (2025), e0324559. https://doi.org/10.1371/journal.pone.0324559

30. M. Abbas, A. Haq, Estimation of finite population distribution function with auxiliary information in a complex survey sampling, *SORT*, **46** (2022), 67–94. https://doi.org/10.2436/20.8080.02.118

31. A. Haq, M. Abbas, M. Khan, Estimation of finite population distribution function in a complex survey sampling, *Commun. Stat. Theor. M.*, **52** (2023), 2574–2596. https://doi.org/10.1080/03610926.2021.1955386

32. N. Nematollahi, M. M. Salehi, R. A. Saba, Two-stage cluster sampling with ranked set sampling in the secondary sampling frame, *Commun. Stat. Theor. M.*, **37** (2008), 2404–2415. https://doi.org/10.1080/03610920801919684

33. M. Abbas, M. Ahmed Shehzad, H. Khurram, M. Rabia, Estimation of the distribution function of a finite population utilizing auxiliary information in the context of non-response within complex survey sampling, *PLoS One*, **20** (2025), e0322660. https://doi.org/10.1371/journal.pone.0322660

34. M. N. Murthy, *Sampling theory and methods*, 1967.

35. A. Haq, J. Brown, E. Moltchanova, Hybrid ranked set sampling scheme, *J. Stat. Comput. Sim.*, **86** (2016), 1–28. https://doi.org/10.1080/00949655.2014.991930

36. M. A. Shehzad, H. Khurram, Z. Iqbal, M. Parveen, M. N. Shabbir, Nutritional status and growth centiles using anthropometric measures of school-aged children and adolescents from multan district, *Arch. Pédiatrie*, **29** (2022), 133–139. https://doi.org/10.1016/j.arcped.2021.11.010

37. M. Yazdi, E. Zarei, S. Adumene, R. Abbassi, P. Rahnamayiezekavat, Uncertainty modeling in risk assessment of digitalized process systems, *Method. Chem. proc. saf.*, **6** (2022), 389–416. https://doi.org/10.1016/bs.mcps.2022.04.005

38. J. Pan, Y. Deng, Y. Yang, Y. Zhang, Location-allocation modelling for rational health planning: Applying a two-step optimization approach to evaluate the spatial accessibility improvement of newly added tertiary hospitals in a metropolitan city of china, *Soc. Sci. Med.*, **338** (2023), 116296. https://doi.org/10.1016/j.socscimed.2023.116296

39. T. A. A. Ali, Z. Xiao, H. Jiang, B. Li, A class of digital integrators based on trigonometric quadrature rules, *IEEE T. Ind. Electron.*, **71** (2024), 6128–6138. https://doi.org/10.1109/TIE.2023.3290247

40. E. Zarei, M. Yazdi, R. Moradi, A. BahooToroody, Expert judgment and uncertainty in sociotechnical systems analysis, In: *Safety causation analysis in sociotechnical systems: Advanced models and techniques*, Cham: Springer, 2024. https://doi.org/10.1007/978-3-031-62470-4_18

41. Y. Lou, M. Cheng, Q. Cao, K. Li, H. Qin, M. Bao, et al., Simultaneous quantification of mirabegron and vibegron in human plasma by hplc-ms/ms and its application in the clinical determination in patients with tumors associated with overactive bladder, *J. Pharmaceut. Biomed.*, **240** (2024), 115937. https://doi.org/10.1016/j.jpba.2023.115937

42. S. Gupta, J. Shabbir, On the use of transformed auxiliary variables in estimating population mean by using two auxiliary variables, *J. Stat. Plan. Infer.*, **137** (2007), 1606–1611. https://doi.org/10.1016/j.jspi.2006.09.008

43. O. Olayiwola, A. Audu, O. Ishaq, S. Olawoore, A. Ibrahim, A class of ratio estimators of a finite population mean using two auxillary variables under two-phase sample scheme, In: *Royal statistical society Nigeria local group annual conference proceedings*, 2020, 80–95.

44. C. Kadilar, H. Cingi, Ratio estimators in stratified random sampling, *Biometrical J.*, **45** (2003), 218–225.

AIMS Press