



## Improving The Accuracy of Area Sampling Frame Estimators Using Unequal Clustered Segment Sampling

Hazanul Zikra

Statistics Indonesia - BPS, Jakarta, Indonesia – hazanul.zikra@bps.go.id

### Abstract

Accurate rice production data is essential for national food security and effective policy planning. The Area Sample Frame (Kerangka Sampel Area – KSA) method has been widely used to estimate rice harvest areas. However, this method has certain limitations, particularly the risk of undercoverage bias when estimating the area across different rice growth stages. This study aims to improve the accuracy of rice area estimations by applying the Unequal Clustered Segment Sampling method as an alternative to the traditional KSA. The proposed method enhances the estimator by excluding non-target segments, that is, spatial points located outside actual rice-growing regions. Using a design-based approach, the estimator accounts for unequal cluster sizes, leading to a more precise representation of field conditions. The results demonstrate that the Unequal Clustered Segment Sampling method significantly reduces bias and improves estimation accuracy compared to the conventional KSA approach. Thus, applying unequal clustered segment sampling designs in KSA-based surveys can yield more reliable and representative estimates, especially in heterogeneous or fragmented agricultural landscapes.

**Keywords:** Area Sampling Frame, Cluster Sampling, Unequal Cluster Size, Rice Production, Agricultural Statistics.

### 1. Introduction

To generate accurate rice production data, the Area Sample Frame (KSA) method has been implemented in Indonesia since 2018 [1]. In estimating the different phases of rice cultivation, Statistics Indonesia (BPS) employs the KSA method by systematically dividing rice fields into sampling units, known as segments, each measuring  $300 \times 300$  meters. These segments are further subdivided into nine equally sized subsegments (equal clusters), each measuring  $100 \times 100$  meters [2].

Although the KSA method relies on objective data collection, further refinement is needed to produce more accurate estimates. The use of square segments may lead bias, as the fixed observation points within these segments can sometimes fall outside the actual paddy field boundaries, leading to inefficiencies and potential frame undercoverage errors [3]. While area frame sampling offers greater protection against non-sampling errors, such as missing or overlapping units in the frame, these types of errors can still occur in area frame surveys, including KSA [4].

The application of Unequal Clustered Sampling with weighting has been shown to be significantly more efficient than its unweighted counterpart. Weighted estimators can adjust for biases resulting from variations in cluster size and the correlation between cluster size and the measured characteristics [5]. In this context, incorporating a correction factor into the estimator, by only points that fall within actual paddy fields when estimating rice growth phases including harvested area, can lead to more efficient and accurate estimates. A cluster refers to a group of observation units that either occur naturally or are purposefully constructed and can function as a single sampling unit [6].



| Type of Estimator  | Methods Code | Type of Estimator                           |
|--|--------------|---|
| Unequal Clustered Segment Sampling with Sample Cluster Weighting     | Method 3     | Consistent Biased from Point of Observation |
| Unequal Clustered Segment Sampling with Population Cluster Weighting | Method 4     | Unbiased from Point of Observation          |

The formulation for calculating the proportion of rice growth phase  $i$  ( $p_i$ ) and its corresponding between-cluster variance ( $s_b^2$ ) under Method 1 is as follow [5]:

$$p_i = \frac{1}{n} \sum_{j=1}^n p_{ij} \quad (1)$$

$$s_b^2 = \frac{1}{(n-1)} \sum_{j=1}^n (p_{ij} - p_i)^2 \quad (2)$$

Where  $p_{ij}$  represents the proportion of phase  $i$  in segment  $j$  and  $n$  is sample size. Meanwhile, the formulation for calculating the proportion of phase  $i$  under Method 2 is shown in equation (3). However, the calculation for the between-cluster variance remains the same as in equation (2).

$$p_i = \frac{1}{n} \sum_{j=1}^n \frac{M_j}{\bar{M}'} p_{ij} = \frac{\sum_{j=1}^n a_{ij}}{\sum_{j=1}^n M_j} \quad (3)$$

Where  $M_j$  is the number of clusters in sample  $j$  and  $\bar{M}'$  is the average number of clusters across all samples. For Method 3, the formulation for calculating the proportion of phase  $i$  is the same as in equation (3). However, the variance for Method 3 is presented in equation (4):

$$s_b^2 = \frac{1}{\bar{M}'^2(n-1)} \sum_{j=1}^n M_j^2 (p_{ij} - p_i)^2 \quad (4)$$

For Method 4, the proportion of phase  $i$  and its corresponding variance are calculated using Equation (5) and (6), as shown below.

$$p_i = \frac{1}{n} \sum_{j=1}^n \frac{M_j}{\bar{M}} p_{ij} \quad (5)$$

$$s_b^2 = \frac{1}{(n-1)} \sum_{j=1}^n \left( \frac{M_j p_{ij}}{\bar{M}} - p_i \right)^2 \quad (6)$$

Where  $\bar{M}$  represents the average number of clusters (subsegments) in the population. The sampling variance, denoted as  $v(p_i)$ , for all methods across all rice growth phases is calculated using Equation (7), as shown below.

$$v(p_i) = \frac{(1-n/N)}{n} s_b^2 \quad (7)$$

### 3. Result and Discussion

#### 3.1. Population of Paddy Field in Denpasar Timur on February 2019

In general, paddy field landscapes in Bali exhibit irregular shapes, are terraced, and tend to be dispersed, and often adjacent to non-agricultural land, a pattern also observed in the East Denpasar District. According to the Ministry of Agrarian Affairs and Spatial Planning (Badan Pertanahan Nasional/BPN), the administrative area of East Denpasar District covers approximately 2,231 hectares [8]. In 2019, the reported extent of paddy fields in the district was around 690 hectares [9]. However, based on the delineation conducted during the Field Work Practice by Politeknik Statistika STIS in February 2019, the identified eligible irrigated rice field area amounted to 479.14 hectares [10]. This presents a discrepancy of 30.56 percent when compared to the paddy field base area reported by the Ministry. The delineation results also show that approximately 21.48 percent of East Denpasar's total area consist of rice fields, indicating notable agricultural potential within the district. This is further supported by economic data: based on the Gross Regional Domestic Product (GRDP) at Current Prices of Denpasar City in 2019, the Agriculture, Forestry, and Fisheries sector ranked fifth in contribution, accounting for 6.48 percent of the total GRDP [11].

**Table 1.** Paddy Field Area by Rice Growth Phase Based on Delineation Results from the Field Work Practice of Politeknik Statistika STIS, Februari 2019 [10].

| Phase                             | Proportion  | Area (Hectare) |
|-----------------------------------|-------------|----------------|
| First Vegetation (V1)             | 0.21        | 101.87         |
| Second Vegetation (V2)            | 0.09        | 41.06          |
| Generative (G)                    | 0.13        | 63.39          |
| Harvested (P)                     | 0.09        | 41.95          |
| Farmland Preparation (PL)         | 0.13        | 62.30          |
| Lodging (B)                       | 0.11        | 51.23          |
| Non-Rice Paddy Field (LL)         | 0.24        | 114.11         |
| Harvest Between Two Surveys (P-2) | ~0          | 2.12           |
| Failed Rice Crop (PS)             | ~0          | 1.11           |
| <b>Total</b>                      | <b>1.00</b> | <b>479.14</b>  |

**Note:** ~0: Has a very small value or approaches zero

The delineation of rice paddy land in February 2019 indicated that the largest proportion of rice cultivation in East Denpasar District was in the non-rice paddy (LL) phase, covering approximately 114.11 hectares, which corresponds to a proportion of 0.24. The harvested area (P) during the same period was around 41.95 hectares, representing roughly 0.09 of the total area. The most extensive standing crop phase was the Early Vegetative (V1) stage, accounting for 101.87 hectares or 21 percent of the paddy fields, followed by the Generative (G) phase with 63.39 hectares (approximately 13 percent). These findings indicate that the majority of rice paddies in East Denpasar were in the early stages of the planting season in February 2019. Conversely, the Potential Crop Failure (PS) phase had the smallest area, covering only 1.11 hectares, representing a proportion close to zero.

In accordance with the established protocols for constructing a square segment sampling frame, a square segment is considered eligible for inclusion in the sampling frame if it contains at least 50 percent of the designated target population [4][12][13]. Applying this criterion to the paddy field population map of East Denpasar resulted in 57 segment units. The proportions and areas corresponding to each rice growth stage, based on observation points falling within the target population (paddy fields), are summarized in Table 2.

**Table 2.** Derived Proportions and Areas of Rice Growth Phases Based on Square Segment Sampling Frame Development.

| Phase                             | Proportion  | Area (Hectare) |
|-----------------------------------|-------------|----------------|
| First Vegetation (V1)             | 0.24        | 116.34         |
| Second Vegetation (V2)            | 0.06        | 29.56          |
| Generative (G)                    | 0.14        | 66.35          |
| Harvested (P)                     | 0.09        | 44.40          |
| Farmland Preparation (PL)         | 0.12        | 55.47          |
| Lodging (B)                       | 0.11        | 52.05          |
| Non-Rice Paddy Field (LL)         | 0.24        | 113.57         |
| Harvest Between Two Surveys (P-2) | ~0          | 1.40           |
| Failed Rice Crop (PS)             | 0.00        | 0.00           |
| <b>Total</b>                      | <b>1.00</b> | <b>479.14</b>  |

**Note:** ~0: Has a very small value or approaches zero

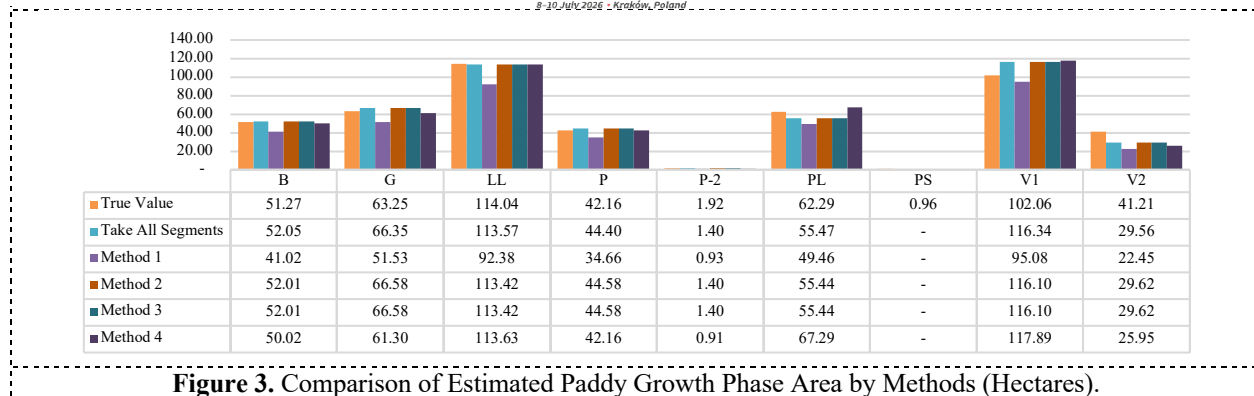
Table 2 illustrates that the parameters derived from the square segment sampling frame differ noticeably from those obtained through direct delineation of rice growth phases without restricting to the target population. This disparity is primarily attributed to undercoverage bias in the KSA Survey, which arises when the sampling frame does not fully capture the target population. Such bias can occur due to segments that fall partially or entirely outside the target area, observation points located outside paddy fields, or the use of uniform square segments that do not conform to the irregular and fragmented distribution of paddy fields, particularly in East Denpasar [14][7].

### 3.2. Comparison of Estimated Rice Growth Phases from Area Frame Survey (KSA)

Based on the simulation results, the average estimated area for rice growth phases using Method 1 (*equal clustered segment sampling*), the method currently implemented by BPS for calculating harvested area, tends to be lower than the true parameter values across all growth phases. This indicates a consistent pattern of underestimation. One contributing factor is the inclusion of observation points that fall outside the actual paddy field target population, which are nevertheless included in the rice growth phase area calculations. This results in a negative bias in the estimation process [7].

In contrast, the application of Method 2 (*unequal clustered segment sampling without weighting*) and Method 3 (*unequal clustered segment sampling with sample weighting*) produced estimates that tended to underestimate several rice growth phases—specifically LL, P-2, PL, PS, and V2—while overestimating others. This pattern is attributed to the “cancel in–cancel out” effect, in which overestimation in certain phases offsets underestimation in others, since the total of all phase proportions is constrained to sum to one.

By comparison, Method 1 does not exhibit this balancing effect; the total estimated proportion does not sum to one due to the inclusion of observation points outside the paddy field target population, which are excluded from rice growth phase calculations, leading to systematic underestimation. Similarly, Method 4 (*unequal clustered segment sampling with population weighting*) also resulted in overestimation for some phases—PL and V1—and underestimation for others—B, G, LL, P-2, PS, and V2. Notably, however, the estimation for the P (harvested) phase under Method 4 was relatively consistent with the actual observed area.



**Figure 3.** Comparison of Estimated Paddy Growth Phase Area by Methods (Hectares).

Simulation sampling results indicate a consistent negative bias in the estimation of rice growth phase area when using Method 1 (equal clustered segment sampling) across all phases. Generally, errors in parameter estimation may arise from two main sources: sampling error and non-sampling error. Sampling error occurs as a natural consequences of selecting a subset (sample) from the population, while non-sampling error arises from factors outside the sampling process. Bias due to non-sampling error can broadly come from specification error, frame error, nonresponse error, measurement error, and processing error [16].

Assuming no errors in observed phases, the shape of rice field polygons, or identification by field officers, the underestimation in Method 1 is predominantly caused by frame under coverage error, which is an error in the construction of the sample frame. Several sources of frame under coverage error that arise when using Method 1, currently employed by Statistics Indonesia-BPS, are as follows: 1) Points falling outside the rice fields (target population), which leads to the estimated area of rice growth phases tending to be underestimated. This occurs because the calculation of rice growth phase proportions still includes these points due to the use of an equal cluster design, 2) The construction of the segment sample frame always places the sub-segment observation points in the center, making it impossible for them to fall within other target populations within the sample sub-segment, and 3) The segment sample frame not covering the entire target population because, to qualify as a segment sample frame, it must contain at least 50 percent rice fields [13][17].

Similar to Method 1, Method 2 (unequal clustered segment sampling without weighting) also produces biased estimations. However, the estimator generated by Method 2 does not include points that fall outside the target population (rice fields). This means the number of sub-segments used in calculating the area of observed phases can vary for each segment. Assuming sub-segments as cluster elements, the segment sample represents a form of unequal cluster size. This is because the number of eligible points within each segment can differ, making the total number of clusters dependent on the number of points that fall within the target population. Compared to Method 1, Method 2 eliminates the bias stemming from points falling outside the target population. The primary sources of bias in Method 2 are: immovable observation points, where observation points are always fixed at the center of the sub-segment (cluster) and static sample frame, the sample frame itself is fixed and does not change.

Meanwhile, Method 3 (unequal clustered segment sampling with sample weighting) produces parameter estimations that are consistently biased. Similar to Method 2, Method 3 estimates rice growth phases by counting points that fall within the target population. However, Method 3 applies a correction factor for the number of elements within each cluster from sample. This means the resulting estimator remains biased, but its bias decreases as the sample size increases, causing the estimator to approach the "take all" value of the observation points from the segment population [18][19]. The sources of bias in Method 3 are immovable observation points, where observation points are always fixed at the center of the sub-segment (cluster), and a static sample frame, meaning the sample frame

itself is fixed and does not change. Consequently, population bias cannot be entirely eliminated. However, the estimator will produce values that approximate the proportion of rice growth phases from all observation points within the segments when samples increased [4].

Method 4 (unequal clustered segment sampling with population weighting) also applies principles similar to Method 3. It excludes points that fall outside the target population, leading to varying sub-segment (cluster) sizes across different segments. Unlike Method 3, Method 4 incorporates a cluster size weight within the population. Estimating the average proportion of rice growth area using Method 4 yields an unbiased estimate of the "take all" proportion of observation points in the segment population, but it requires information on the average number of clusters from the segment population [5][6].

### 3.3. Comparison of Biased of Estimation Area Growth Phase from Area Frame Survey (KSA)

The average bias across all sampling simulations indicates that Method 1 (equal clustered segment sampling) generally results in underestimation for all phases [13]. The most significant bias was observed in the LL phase, reaching -21.65 hectares, followed by the V2 phase, which experienced an under coverage bias of -18.75 hectares. The absolute bias for Method 1 is 10.18 hectares. As discussed earlier, Method 1 indeed produces an estimator for rice growth phase area that is biased under coverage.

Conversely, using Method 2 (unequal clustered segment sampling without weighting) and Method 3 (unequal clustered segment sampling with sample weighting) resulted in generally smaller biases. The most significant under coverage bias for both Method 2 and Method 3 was observed in the V2 phase, at -11.59 hectares. Meanwhile, the highest over coverage bias for both methods occurred in the V1 phase, with a bias level reaching 14.04 hectares. The biases produced by Method 2 and Method 3 are identical because their distinction lies solely in the calculation of sampling variance and Mean Squared Error (MSE). However, the estimation of rice growth phase area is relatively the same for both methods, as they only consider observation points that fall within the target population, i.e., paddy fields. The absolute bias for Method 2 and Method 3 are 4.56 hectares.

The average absolute bias for Method 4 is 4.63 hectares. The highest estimator bias was observed in the V1 phase, with a value of 15.83 hectares, while the deepest bias was in the V2 area estimator, at -15.26 hectares. Theoretically, Method 4's application should yield an unbiased estimator for the "take all" observation points. However, any resulting bias is potentially due to under coverage of the estimator when compared to the total area of all rice field polygons. The bias observed across different observation phases also exhibits a "cancel in, cancel out" effect. This happens because the observation points considered in calculating the area of rice growth phases only include those that fall within the rice fields [19].

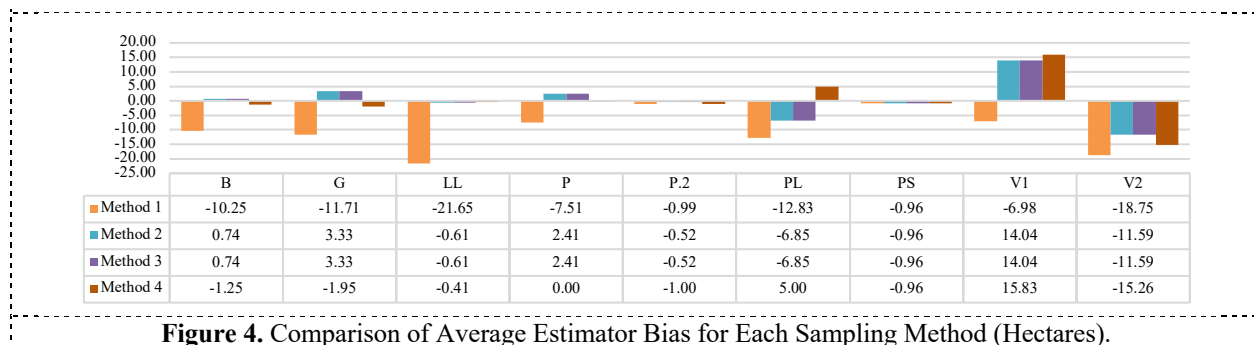


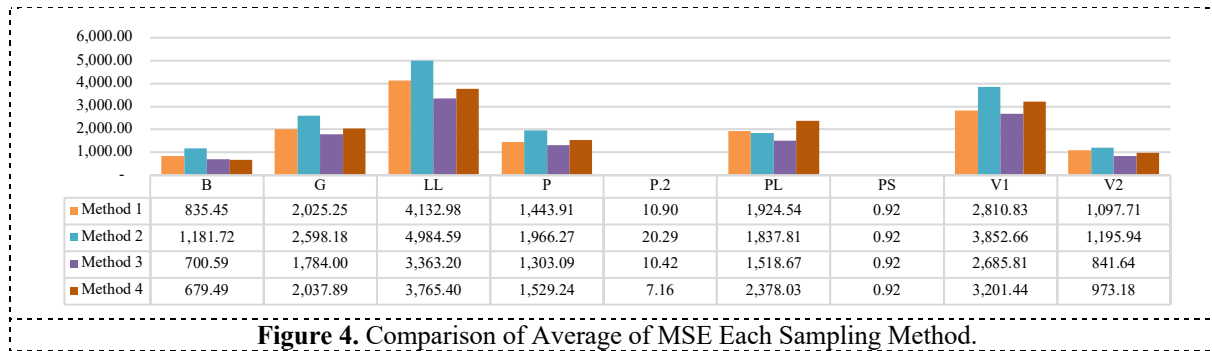
Figure 4. Comparison of Average Estimator Bias for Each Sampling Method (Hectares).

### 3.4. Comparison of Mean Square Error (MSE) Growth Phase from Area Frame Survey (KSA)

To compare the accuracy across different sampling designs, the Mean Squared Error (MSE) is utilized. MSE measures the difference between an estimator and the parameter being estimated. Mathematically, MSE is calculated using the following formula [16]:

$$MSE = Bias^2 + Variance \quad (8)$$

Based on that formulation, the MSE of all sampling simulations is shown in the following figure.



Based on the MSE calculations for each method at every rice growth stage, it is evident that applying Method 2 results in a high MSE. This is triggered by the variability in the proportion of rice growth stages, which stems from applying an unequal clustered size to each segment without a correction factor for the number of cluster elements (subsegments).

Meanwhile, Method 3 yields an average MSE value that is lower than the other methods. This is because variability leads to a lower sampling variance due to a correction factor for the differences in the number of clusters in each segment, thus resulting in a lower sampling variance.

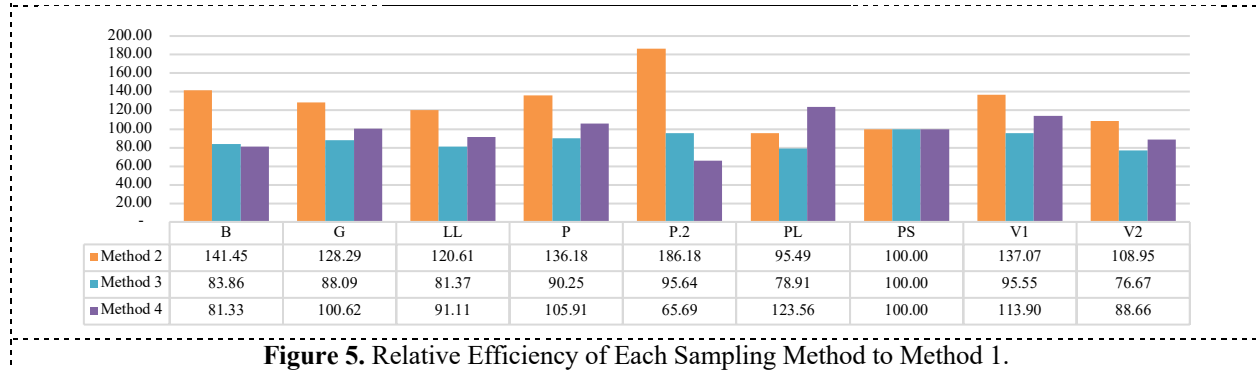
The application of Method 4 yields a relatively higher MSE compared to the other methods, as the estimated proportions generated are relatively diverse even though a weighting factor for cluster elements within the population has been applied. Theoretically, the variance of an average estimator with population weights is often greater than the variance with sample weights, especially when the cluster sizes vary significantly and there is a correlation between cluster size and the average cluster value. [5]. For Method 1, the elevated Mean Squared Error (MSE) is attributed not only to sampling variance but also to a consistent underestimation of the parameter's true value. While this leads to a higher relative MSE, it remains considerably lower than the MSE observed in Method 2, which is a result of employing an equal cluster size design.

### 3.5. Relative Efficiency

To determine the efficiency of the estimation results for each sampling method relative to Method 1 or the Statistics Indonesia-BPS KSA method, relative efficiency (RE) is calculated by comparing the Mean Squared Error (MSE) of each KSA sampling method with the MSE of the Method 1. The optimal sampling option is determined by the lowest average relative efficiency (RE). Based on the simulation results, the average relative efficiency across all rice growth phases using Method 2 (unequal clustered segment sampling without weighting) is greater than 100. This indicates that the estimator is not more efficient compared to the current Statistics Indonesia-BPS KSA method. The average RE for Method 2 is 128.25.

**Table 3.** Average Relative Efficiency Across All Phases.

| Methods  | Average of RE |
|----------|---------------|
| Method 2 | 128.25        |
| Method 3 | 87.82         |
| Method 4 | 96.75         |



**Figure 5.** Relative Efficiency of Each Sampling Method to Method 1.

In contrast to Method 2 (unequal clustered segment sampling without weighting), Method 3 (unequal clustered segment sampling with sample weighting) yielded an average Relative Efficiency (RE) of 87.82. Generally, the RE values for all phases using Method 3 were less than 100. This indicates that, overall, Method 3 offers significantly better efficiency compared to the current KSA method employed by Statistics Indonesia-BPS.

Similar to Method 3, Method 4 (unequal clustered segment sampling with population weighting) also offers improved estimator efficiency compared to Method 1 (equal clustered segment sampling). This is evident from its average Relative Efficiency of less than 100, and it is generally even lower than Method 3, with an overall average of 96.75.

When considering the potential application of methods for estimating rice growth phase area within Statistics Indonesia-BPS, Method 3 (unequal clustered segment sampling with sample weighting) is more suitable. This is because it only requires weighting information from the sample itself [20]. Although Method 4 (unequal clustered segment sampling with population weighting) offers better efficiency, its implementation for estimating harvest area is less practical due to the general unavailability of population segment information. To effectively apply Method 4, the exact number of observation points in the population must be known to achieve a higher level of estimator precision.

#### 4. Conclusion

Based on the simulation results, the application of Method 3 (unequal clustered segment sampling with sample weighting) offers a relatively more efficient estimation of rice growth phase areas compared to the current KSA Method used by BPS. This is supported by a Relative efficiency (RE) value of less than 100, indicating improved precision over the traditional method.

Moreover, the use of unequal clustered segment sampling with sample weighting is particularly effective in reducing frame undercoverage bias, a key limitation observed in KSA survey estimations. The incorporation of a correction factor based on cluster size allows this method to adjust for uneven distributions of observation points and is practically feasible for implementation by Statistics Indonesia (BPS), especially given the limited availability of



complete segment population data across all regions. Scientifically, this approach provides estimates with a lower error rate than equal cluster sampling, making it a viable alternative for improving the accuracy of rice area statistics.

However, despite the rigorous simulation design, this research is subject to several limitations. The data utilized are specific to East Denpasar, Bali, and may not reflect the characteristics of rice cultivation in other regions or in non-paddy field environments. Therefore, further validation through expanded simulations and field testing in diverse rice-producing regions is recommended to ensure the method's generalizability and robustness.

## References

- [1] Badan Pusat Statistik, *Pedoman Teknis Pendataan Statistik Pertanian Tanaman Pangan Terintegrasi dengan Metode Kerangka Sampel Area*. Jakarta: BPS, 2018.
- [2] Badan Pusat Statistik, *Upaya Perbaikan Data Padi dengan Metode Kerangka Sampel Area (KSA)*. Jakarta: BPS, 2018.
- [3] J. Gallego, *Area sampling frames for agricultural and environmental surveys*. Publications Office, 2015.
- [4] F. J. Gallego and J. Delincé, *The European Land Use and Cover Area-frame Statistical Survey*. New York: Wiley, 2010.
- [5] W. G. . Cochran, *Sampling techniques*. Wiley, 2005.
- [6] A. Asra and A. Prasetyo, *Pengambilan Sampel dalam Penelitian Survei*. Jakarta: PT RajaGrafindo Persada, 2015.
- [7] H. Zikra and W. P. Buana, "Analisis Perbandingan Desain Sampling Survei Kerangka Sampel Area (KSA)," in *Seminar Nasional Official Statistics*, Jakarta: Politeknik Statistika STIS, Sep. 2020, pp. 1326–1336.
- [8] BPS Kota Denpasar, *Kota Denpasar Dalam Angka Denpasar 2020*. Denpasar: BPS Kota Denpasar, 2020.
- [9] BPS Kota Denpasar, *Kecamatan Denpasar Timur dalam Angka 2020*. BPS Kota Denpasar, 2020.
- [10] PKL STIS 58, *Studi Akurasi KSA dengan Pendekatan Deliniasi*. Jakarta: Politeknik Statistika STIS, 2019.
- [11] Badan Pusat Statistik Kota Denpasar, *Produk Domestik Regional Bruto Kota Denpasar Menurut Lapangan Usaha 2018-2022*. 2025.
- [12] Mubekti and L. Sumargana, "Pendekatan Kerangka Sampel Area untuk Estimasi dan Peramalan Produksi Padi," *Pangan*, vol. 25, no. 2, pp. 71–146, 2016.
- [13] F. J. . Gallego, *Sampling frames of square segments : An agricultural information system for the European Union*. OPOCE : European Commission, 1995.
- [14] H. Zikra and W. P. Buana, "Penggunaan Systematic Point Sample Sebagai Area Master Frame dalam Mengestimasi Luas Panen Padi (Studi Simulasi Sampling di Kecamatan Denpasar Timur Tahun 2019)," *J. Stat. dan Apl.*, vol. 6, no. 1, 2022.
- [15] Q. Dong, J. Liu, L. Wang, Z. Chen, and J. Gallego, "Estimating crop area at county level on the North China plain with an indirect sampling of segments and an adapted regression estimator," *Sensors (Switzerland)*, vol. 17, no. 11, 2017, doi: 10.3390/s17112638.
- [16] P. P. Biemer and L. E. Lyberg, "Introduction to Survey Quality," 2003. [Online]. Available: [www.copyright.com](http://www.copyright.com).
- [17] F. Javier, G. Pinilla, and E. Commission, *Sampling Frames of Square*, no. January 1995. 2017.
- [18] S. L. Lohr, *Sampling: Design and Analysis*. Brooks/Cole, 2010.
- [19] L. Kish, *Survey Dampling*. New York: Wiley-Interscience, 1965.
- [20] J. Delince, F. Javier, G. Pinilla, and E. Commission, "Technical Report Series GO-39-2017 Agricultural Master Sampling Frames in Practice Lessons learned from international field experiments and case studies," no. January 2019, 2018, doi: 10.13140/RG.2.2.31374.20808.