

# Misclassification in stratified agricultural area samples

Cristiano Ferraz<sup>1</sup>, Raydonal Ospina<sup>1\*</sup>, André Leite<sup>1</sup> and Hemílio Coelho<sup>1,2</sup>

<sup>1</sup>CAST - Computational Agricultural Statistics Laboratory (CAST),  
Department of Statistics, Universidade Federal de Pernambuco, Recife, Brazil

<sup>2</sup>Department of Statistics, Universidade Federal da Paraíba, Paraíba, Brazil

---

## Abstract

Stratification is a sampling technique widely used in surveys to improve efficiency of estimators. Strata building processes are subject to different types of errors that can lead to misclassification of sampled units, a discrepancy between the stratum from which a unit was selected, and the stratum the unit belongs to in reality. This paper investigates the problem motivated by surveys using stratified area frames of square segments to generate agricultural statistics. Estimators coping with the problem are introduced and their statistical performance is investigated using a Monte Carlo simulation experiment. The study rely on a real-case motivated scenario in which area frames of square segments were applied to surveys carried out in two Brazilian municipalities aiming at comparing different sampling design strategies to generate efficient agricultural statistics. Simulation results indicate that the adoption of a naive estimator can introduce bias and inflate variance. It also indicates that in the absence of the needed auxiliary information to use a post-stratified estimator, the best choice is to keep the design-based original sample stratification estimator.

**Keywords and Phrases:** errors-in-stratification, agricultural surveys, master frame, area sampling, area frame of square segments, crowd-sourcing, post-stratification

---

## 1 Introduction

Misclassification problems related to survey stratification occur when field collected data differs from auxiliary information used to stratify the population. In practice this means sampled units present characteristics that classify them in strata that do not correspond to the ones from which they were originally selected. When coping with strata misclassification in practice, options considered include keeping the misclassified sampled units in the original stratum or move them to the actual one. In addition, the choice of an appropriate estimator for the parameter of interest needs to take into account the type of information available.

Several situations can lead to misclassification problems in practice. Mulrow and Woodburn (1990) describe a study of the effect of errors in stratification in a problem of business taxation. In their case, a small simulation study was carried out investigating bias effects on estimates, but not addressing variance problems. Lamas et al. (2010) and Abreu et al. (2010) report an example of misclassification of tracts as agricultural or non-agricultural in the June Area Survey, carried out by NASS, the US National Agricultural Statistics Service, leading to discrepancies between the estimated number of agricultural farms using the survey data, and the agricultural census at that time. A generalized linear model was used to model under-counting and to

---

\*Correspondence: Raydonal Ospina. E-mail: [raydonal@de.ufpe.br](mailto:raydonal@de.ufpe.br)

provide corrections for estimates. Jang et al. (2009) describe misclassification in stratification exemplifying on the discrepancies observed at the National Survey of Recent College Graduates (NSRCG), between the race/ethnicity registered in administrative records, and the self reported one. Another general situation leading to strata misclassification problems happens when there is a large time lag between the period the sample design was developed, the sample was selected, and the period of the field work for data collection. Sometimes this time lag is related to the stratification process itself. Agricultural surveys using satellite imagery to foster stratification of segments of lands, for example, are subject to this source of error when the image used for strata classification is not updated. This can represent in practice a time lag discrepancy from the information used to build the strata and the reality found in the field. In this paper, misclassification is studied in the context of agricultural surveys based on stratified area samples of squared segments. The study rely on a real-case motivated scenario where the efficiency of area frames of square segments, as master sampling frames for agricultural statistics, was investigated by surveys carried out in two Brazilian municipalities: Goiana, and Santos Dumont. The surveys were part of the studies carried out by the Food and Agriculture Organization of the United Nations - FAO's Global Strategy to Improve Agricultural and Rural Statistics (GSARS), in 2016 and 2017, described in detail in Ferraz et al. (2018). The stratification process used supervised classification of image points within squared segments, relying on the idea of crowd-sourcing to foster stratification. Discrepancies between the stratification generated by imagery analysis and the stratification based on field collected data define the misclassification problem that motivates this paper. Several estimators that could cope with the situation are introduced and their performances analysed via Monte Carlo simulation. The computational experiment replicates stratified samples drawn from an artificial population built to resemble major characteristics of Goiana, simulating a misclassification process with the same rates found in practice, by data collection on the field.

## 2 Misclassification in crowd-sourcing

One of many aspects investigated by the GSARS surveys carried out in Brazil, and Nepal was stratification efficiency of area frames of square segments based on free satellite imagery, and supervised classification using crowd-sourcing. Area frame stratification applied to Brazil's experiment used free imagery resources from Google Earth and Open Street Map (powered by [esri.com](http://esri.com)) to support photo interpretation of points in square segments, by volunteers. The major advantage of this method relies on its low cost of implementation in terms of budget and timing to achieve stratification of the full territory. According to IBGE ([cidades.ibge.gov.br](http://cidades.ibge.gov.br)), Goiana has an area of 445.814 Km<sup>2</sup>, and Santos Dumont, 637.373 Km<sup>2</sup>. Area frames for Goiana and Santos Dumont used square segments of 49 hectares, and within each segment, a sample of 5 points was taken, following the pattern shown in Figure 1(a). The land cover in each point was assessed, by volunteers, and classified according to two classes: "*Cropland*" or "*Non-Cropland*". Figure 1(b). shows an example of classification screen used in the experiment. After classifying all the points, square segments were grouped into strata according to Table 1. Using limited resources of six volunteers, stratification of Goiana was achieved in one week, while the same task for Santos Dumont, a larger territory, was completed in three weeks, using only three volunteers. Although the number of people involved in assessing images for the experiment is extremely low, the idea of crowd-sourcing, were a much larger number of people can be involved in the process, is the motivation to make the method feasible for much larger territories, justifying the use of the term even for the small GSARS study.

The process used to stratify the grid of square segments, although of low cost, is subjected to at least three major sources of errors: the degree of experience of the volunteer photo interpreter, the image quality and the time the image was acquired. The Photo interpretation was carried with a rough rate of 500 points per day approximately, 50,000 per week, with six volunteers.

Table 1: Strata definitions

Strata	Description
1. Highly cropland	Segments with 4 or 5 points classified as “Cropland.”
2. Cropland	Segments with 2 or 3 points classified as “Cropland.”
3. Non-cropland	Segments with at most 1 point classified as “Cropland.”

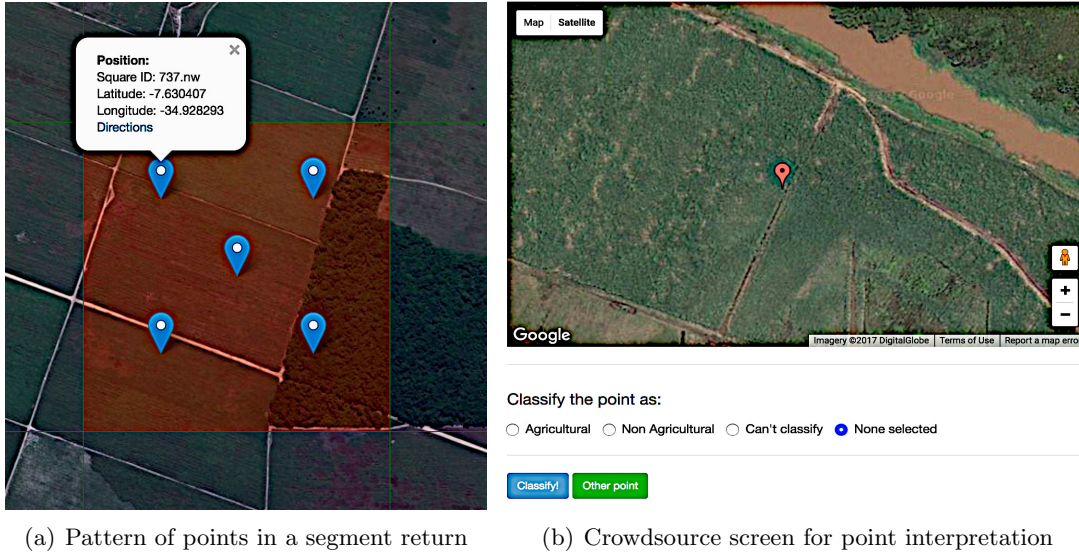


Figure 1: Classification of points in segments

Images from 2010 and 2016 were used to stratify Goiana and Santos Dumont, respectively. Although it is possible to use even more up to date images than in Santos Dumont (based on European satellite Sentinel-2, for instance) this would not necessarily mean misclassification problems would be reduced to the point they would not occur. Therefore, to investigate the potential impact of strata misclassification is necessary and relevant.

Area frames of square segments on both municipalities were stratified and a sample of 60 segments in each city was allocated to the strata according to Table 2.

Table 2: Original square segments classification in strata by the time the sample was selected

Strata	Municipality			
	Goiana		Santos Dumont	
	Frame	Sample	Frame	Sample
1. Highly Cropland	314	42	31	29
2. Cropland	232	15	288	16
3. Non-Cropland	376	3	1078	15

Tables 3 and 4 introduce sample rates of errors in strata formation. These numbers provide estimates of the proportions of segments misclassified in the whole grid of square segments. Based on them, approximately 58% of the segments in Goiana’s area frame were classified in the right strata, while 79% of the segments in Santos Dumont’s area frame were correctly classified. Using data information from both municipalities, the method classified 71% of the segments in the correct strata. These data are used to study options to correct estimates taking into account the fact that the original stratification process was not 100% accurate.

Table 3: Goiana strata misclassification analysis  
Satellite (lines) by Field (columns)

<b>Count</b>				
Total %		<b>Highly Cropland</b>	<b>Cropland</b>	<b>Non-Cropland</b>
Col %		<b>(1)</b>	<b>(2)</b>	<b>(3)</b>
Row %				<b>Total</b>
		<b>27</b>	<b>5</b>	<b>10</b>
<b>Highly Cropland</b>		45.00	8.33	16.67
<b>(1)</b>		87.10	45.45	55.56
		64.29	11.90	23.81
		<b>4</b>	<b>5</b>	<b>6</b>
<b>Cropland</b>		6.67	8.33	10.00
<b>(2)</b>		12.90	45.45	33.33
		26.67	33.33	40.00
		<b>0</b>	<b>1</b>	<b>2</b>
<b>Non-Cropland</b>		0.00	1.67	3.33
<b>(3)</b>		0.00	9.09	11.11
		0.00	33.33	66.67
		31	11	18
<b>Total</b>		51.67	18.33	30.00
				<b>60</b>

Let  $p_{hj}$  be the sample proportion of segments classified in stratum  $j$  in the field, given all the segments classified in stratum  $h$  when using satellite imagery photo interpretation. These numbers correspond to the row percent values in Tables 3 and 4. In Goiana, for instance,  $p_{12} = 0.119$  is the proportion of segments in the sample that were originally selected from stratum 1 (*Highly cropland*) but after data collection, were classified in stratum 2 (*Cropland*). This same proportion for Santos Dumont was 0.2759. It is noted that, in Goiana,  $p_{31} = 0$ , and in Santos Dumont,  $p_{21} = p_{32} = p_{31} = 0$ .

### 3 Correcting estimates

Let  $t_c$  be the total area cultivated with a given crop  $c$ , in a given municipality.  $t_c$  can be expressed, based on strata defined over the population, as

$$t_c = \sum_{h=1}^H \sum_{k \in U_h} y_k,$$

where  $h$  is the index identifying the original stratum (based on photo interpretation),  $U_h$  is the set of all segments of the population from stratum  $h$  and  $y_k$  is the area cultivated with crop  $c$

Table 4: Santos Dumont strata misclassification analysis  
Satellite (lines) by Field (columns)

<b>Count</b>				
Total %		<b>Highly Cropland</b>	<b>Cropland</b>	<b>Non-Cropland</b>
Col %		<b>(1)</b>	<b>(2)</b>	<b>(3)</b>
Row %				<b>Total</b>
	<b>7</b>	<b>8</b>	<b>14</b>	
<b>Highly Cropland</b>	11.67	13.33	23.33	29
<b>(1)</b>	100.00	88.89	31.80	48.33
	24.14	27.59	48.27	
	<b>0</b>	<b>1</b>	<b>15</b>	
<b>Cropland</b>	0.00	1.66	25.00	16
<b>(2)</b>	0.00	11.11	34.10	26,67
	0.00	6.25	93.75	
	<b>0</b>	<b>0</b>	<b>15</b>	
<b>Non-Cropland</b>	0.00	0.00	25.00	15
<b>(3)</b>	0.00	0.00	34.10	25.00
	0.00	0.00	100.00	
<b>Total</b>	7	9	44	60
	16.67	15.00	73.33	

in segment  $k \in U_h$ . The general class of homogeneous linear estimators for this parameter is composed by estimators that can be written as

$$\hat{t}_c = \sum_{h=1}^H \sum_{k \in U_h} I_k w_k y_k$$

where  $I_k$  is the sample inclusion indicator for element  $k$ , and  $w_k$  does not contain any information concerning the response variable. In such conditions, sampling is non-informative and the mean square error of  $\hat{t}_c$  can be written as

$$\hat{t}_c = \sum_{h=1}^H \sum_{g=1}^H \sum_{k \in U_h} \sum_{l \in U_g} y_k y_l E_p(I_k w_k - 1)(I_l w_l - 1),$$

where  $E_p(\cdot)$  is the expectation with respect to the sample design  $p(\cdot)$ . When  $w_k = \pi_k^{-1}$ , with  $\pi_k$  as the inclusion probability of element  $k$ , the estimator  $\hat{t}_c$  corresponds to the traditional Horvitz-Thompson estimator.

Consider  $N_{hj}$  as the number of segments in the area frame that were originally classified in stratum  $h$ , based on image interpretation, and later classified in stratum  $j$ , based on field observation. Table 5–(a) shows the crosstabulation of  $N_{hj}$  illustrating the relationship between the misclassified cells and the actual strata totals  $N_{+j}$ . Let  $n_{hj}$  be the number of segments in the area frame sample that were originally classified in stratum  $h$ , based on image interpretation, and later classified in stratum  $j$ , based on field observation. Table 5–(b) shows the crosstabulation of  $n_{jh}$  illustrating the relationship between the misclassified sample cells and the actual marginal totals  $n_{+j}$ . If no misclassification occurs, then the design strata is the same as the actual strata, with  $N_{h+} = N_{+j}$  and  $n_{h+} = n_{+j}$ .

Table 5: Crosstabulation of stratum sizes

Design Strata	Actual Strata				Design Strata	Actual Strata			
	1	2	3	Total		1	2	3	Total
1	$N_{11}$	$N_{12}$	$N_{13}$	$N_{1+}$	1	$n_{11}$	$n_{12}$	$n_{13}$	$n_{1+}$
2	$N_{21}$	$N_{22}$	$N_{23}$	$N_{2+}$	2	$n_{21}$	$n_{22}$	$n_{23}$	$n_{2+}$
3	$N_{31}$	$N_{32}$	$N_{33}$	$N_{3+}$	3	$n_{31}$	$n_{32}$	$n_{33}$	$n_{3+}$
Total	$N_{+1}$	$N_{+2}$	$N_{+3}$	$N_{++}$	Total	$n_{+1}$	$n_{+2}$	$n_{+3}$	$n_{++}$

(a) Population

(b) Sample

In practice, sample selection is done within strata formed by the rows of Table 5–(a), and when misclassification is present, the values of  $N_{hj}$  for  $h$  differing to  $j$ , are not zero. Table 3 for example, shows that in Goiana,  $N_{12} = 5$ , meaning 5 squares originally selected from stratum 1 were actually found in stratum 2 in the field. In such cases the population can be partitioned in two groups. One based on the design strata, and another based on the field strata. Taking Goiana again as an example, Table 3 shows that according to the design strata, the population of segments is partitioned in groups of size  $N_1 = N_{1+} = 42$ ,  $N_2 = N_{2+} = 15$  and  $N_3 = N_{3+} = 42$ , while according to the field strata, the population is partitioned in groups of size  $N_1 = N_{+1} = 31$ ,  $N_2 = N_{+2} = 11$  and  $N_3 = N_{+3} = 18$ . In this situation, the sample set can also be partitioned into two groups of strata, corresponding to the design and the field information. Let the set of sampled units classified by the field strata, of size  $n_{+j}$ , be denoted by  $A_j$ , for  $j = 1, \dots, H$  and keep the notation  $S_h$  for the set of sampled units of size  $n_{h+}$ , classified by the design strata, with  $h = 1, \dots, H$ . Estimators coping with errors in stratification consider either using  $A_j$  or  $S_h$  as the sample of units, and may modify the form of the sampling weight  $w_k$  as well. This paper considers the same set of estimators described in Mulrow and Woodburn (1990), this time related to the use of an area sampling frame design and strata errors observed in crowd-sourcing based classification. The estimators' description, based on a simple random sampling design, are presented next.

1. **Basic estimator:** The basic estimator corresponds to ignore errors in stratification, keeping the design strata, and choosing to use  $w_k = N_{h+}/n_{h+}$ , for  $k \in S_h$ . Thus, this estimator can be written as follows:

$$\hat{t}_c = \sum_{h=1}^H \sum_{k \in S_h} (N_{h+}/n_{h+}) y_k = \sum_{h=1}^H N_{h+} \bar{y}_{h+},$$

where  $\bar{y}_{h+}$  is the mean of all the segments in the sample originally selected from stratum  $h$ . In this case, no matter the sampling units actually belong to another stratum, they are kept in the design stratum anyway.

2. **Unweighted estimator:** The unweighted estimator corresponds to proceed corrections to the basic estimator using only the field stratum observation in the sample. According to Table 3, in Goiana, a sample of size  $n_{1+} = 42$  was selected from stratum 1, and from these, only  $n_{+1} = 31$  remained in stratum 1, according to the field information. Stratum 1, by design, was composed by  $N_{1+} = 314$  segments (see Table 2). Then for a segment  $k$  in

stratum 1,  $w_k = [314 - (42 - 31)]/31$ . Thus,  $w_k = [N_{j+} - (n_{j+} - n_{+j})]/n_{+j}$ , for  $k \in A_j$ . Note that corrections are made on both, the stratum population and the respective sample size. These corrections, however, take into account only those observed sampled units that changed strata. No tentative is made to proceed corrections to the non-sampled segment numbers. Thus, this estimator can be written as follows:

$$\hat{t}_c = \sum_{j=1}^H \sum_{k \in A_j} [N_{j+} - (n_{j+} - n_{+j})] y_k / n_{+j} = \sum_{j=1}^H [N_{j+} - (n_{j+} - n_{+j})] \bar{y}_{+j},$$

where  $\bar{y}_{+j}$  is the mean of the segments in  $A_j$ .

3. **Weighted estimator:** The weighted estimator uses sample field data to correct information for the observed and the non-observed segments. Let  $\hat{N}_{+j}$  be an estimator for the number of segments over the whole area frame that belong to stratum  $j$ . Let  $w_k = \hat{N}_{+j}/n_{+j}$ , for  $k \in A_j$ . Then,

$$\hat{t}_c = \sum_{j=1}^H \sum_{k \in A_j} \left( \hat{N}_{+j} / n_{+j} \right) y_k = \sum_{j=1}^H \hat{N}_{+j} \bar{y}_{+j}$$

where  $\hat{N}_{+j} = \sum_{h=1}^H N_{jh} p_{hj}$  and  $p_{hj}$  are the proportions of segments in the sample that were originally classified in stratum  $h$  but actually belong to stratum  $j$ .

4. **Post-Stratified estimator:** The post-stratified estimator uses the most updated information for both, the sample and the population sizes. So, for a segment  $k$  in  $A_j$ ,  $w_k = N_{+j}/n_{+j}$ . In practice, no information about the true values  $N_{+j}$  is available, but this estimator is important as a benchmark for the performances of the remaining ones.

## 4 Models and Simulation

In order to evaluate the statistical performances of the considered estimators, an artificial population resembling the field characteristics of Goiana was built to provide support for a Monte Carlo simulation. The population was covered by an area frame of square segments with 922 segments, keeping the strata population sizes presented in Table 2: 314, 232 e 376, for strata 1, 2 and 3, respectively. Several scenarios, replicating the same rates of misclassification found in Goiana's experiment (see Table 2), were considered. FAO's GSARS' experiments used a sample of size 60 due to budget constraints. In this study, the sample sizes investigated were of 120, 180 and 240. In addition, two stratum allocation rules were considered: uniform and proportional to the strata sizes. The estimators described in Section 3 were applied to each investigating scenario and their distribution studied based on 10,000 Monte Carlo replications. The observations and respective misclassification cases were generated according to the models described next. Population means for the area cultivated with sugarcane per square in each stratum were obtained empirically from Goiana's study:  $\mu_1 = 28.35$ ,  $\mu_2 = 22.81$  and  $\mu_3 = 1.42$ .

### 4.1 Parameter's design

Let  $y_k$  be the area cultivated with sugarcane observed in segment  $k$ . Define  $\delta_{hj}^k$  as an indicator variable that segment  $k$  was sampled from the design stratum  $h$ , and later reclassified, in the field stratum  $j$ , so that  $\delta_{hj}^k \sim \text{Bernoulli}(p_{hj})$ . Then the observations and respective errors in strata were generated according to the following models, with  $\varepsilon_k \sim \mathcal{N}(0, 1)$ :

- **Stratum 1 Model - Highly-Cropland:**

314 observations were generating according to:

$$\xi_1 : y_k = \mu_1 - (\mu_1 - \mu_2) \delta_{12}^k - (\mu_1 - \mu_3) \delta_{13}^k + \varepsilon_k,$$

- **Stratum 2 Model - Cropland:**

232 observations were generating according to:

$$\xi_2 : y_k = \mu_2 + (\mu_1 - \mu_2) \delta_{21}^k - (\mu_2 - \mu_3) \delta_{23}^k + \varepsilon_k,$$

- **Stratum 3 Model - Non-Cropland:**

376 observations were generating according to:

$$\xi_3 : y_k = \mu_3 + (\mu_2 - \mu_3) \delta_{32}^k + (\mu_1 - \mu_3) \delta_{31}^k + \varepsilon_k,$$

The parameter values generated according to models  $\xi_1$ ,  $\xi_2$ , and  $\xi_3$  are: 4993.65 (Stratum I), 3718.92 (Stratum II), and 1345.46 (Stratum III), leading to a population total of  $t_c = 10,058.03$  ha.

## 5 Simulation results

The summary statistics for the results of the Monte Carlo experiment is presented in

Table 6: Summary Statistics of 10000 Monte Carlo Replications. Sample size: 60

Stratum	Estimator	Estimador	Bias	Rel.bias	EQM	SD
I	Basic	5027.53	33.88	0.01	364673.14	602.96
	Unweigthed	6796.53	1802.88	0.36	3616724.91	605.30
	Weigthed	6452.70	1459.05	0.29	4947341.77	1678.92
	PostStrat	6462.59	1468.94	0.29	2487798.04	574.49
II	Basic	3728.02	9.10	0.00	460885.40	678.86
	Unweigthed	3447.14	-271.78	-0.07	461695.72	622.79
	Weigthed	4716.95	998.03	0.27	3720555.89	1650.69
	PostStrat	4782.97	1064.05	0.29	1878180.97	863.74
III	Basic	1333.91	-11.55	-0.01	3231841.80	1797.79
	Unweigthed	501.41	-844.05	-0.63	725494.98	114.35
	Weigthed	400.83	-944.63	-0.70	917662.48	159.16
	PostStrat	403.83	-941.63	-0.70	895157.49	92.12
Total	Basic	10089.46	31.43	0.00	3993771.56	1998.30
	Unweigthed	10745.09	687.05	0.07	1217802.85	863.62
	Weigthed	11570.48	1512.45	0.15	6098636.32	1952.31
	PostStrat	11649.39	1591.36	0.16	3597847.92	1032.24



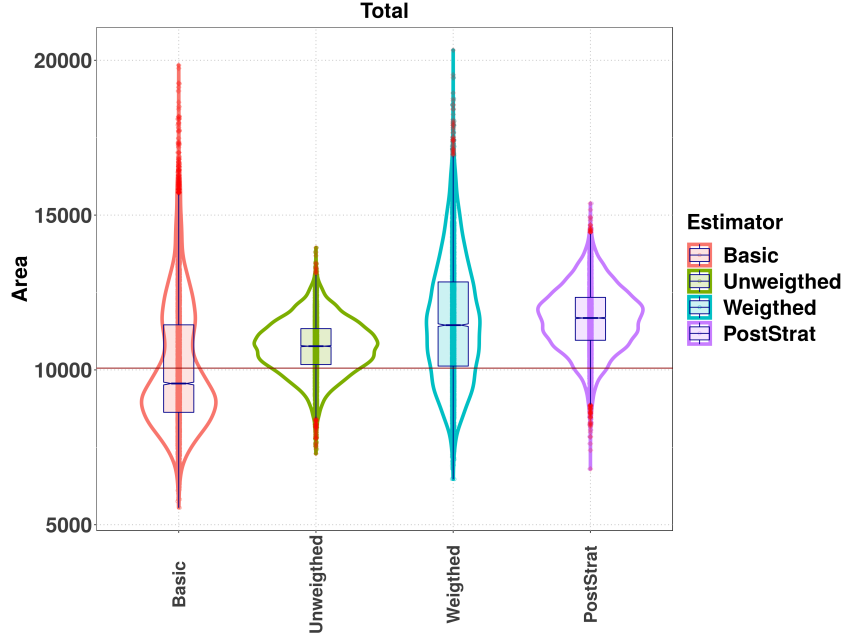


Figure 2: Violin-plots with Boxplots of Estimators Performances for 10000 Monte Carlo Replicate

Table 7: Summary Statistics of 10000 Monte Carlo Replications. Sample size: 120

Stratum	Estimator	Estimador	Bias	Rel.bias	EQM	SD
I	Basic	5025.72	32.07	0.01	366861.66	604.87
	Unweighed	6478.91	1485.26	0.30	2540863.88	578.71
	Weigthed	5995.75	1002.09	0.20	1659163.70	809.34
	PostStrat	6000.43	1006.78	0.20	1301617.28	536.69
II	Basic	3722.26	3.34	0.00	210164.10	458.45
	Unweighed	2938.94	-779.98	-0.21	765804.78	396.81
	Weigthed	3893.75	174.83	0.05	428121.39	630.55
	PostStrat	3895.87	176.95	0.05	308872.35	526.87
III	Basic	1346.26	0.80	0.00	180997.58	425.46
	Unweighed	280.66	-1064.80	-0.79	1136570.26	52.69
	Weigthed	237.93	-1107.53	-0.82	1229356.33	52.27
	PostStrat	237.66	-1107.80	-0.82	1229193.49	44.37
Total	Basic	10094.24	36.21	0.00	749496.33	865.02
	Unweighed	9698.51	-359.52	-0.04	599497.28	685.78
	Weigthed	10127.42	69.39	0.01	753101.96	865.08
	PostStrat	10133.96	75.93	0.01	552531.14	739.47

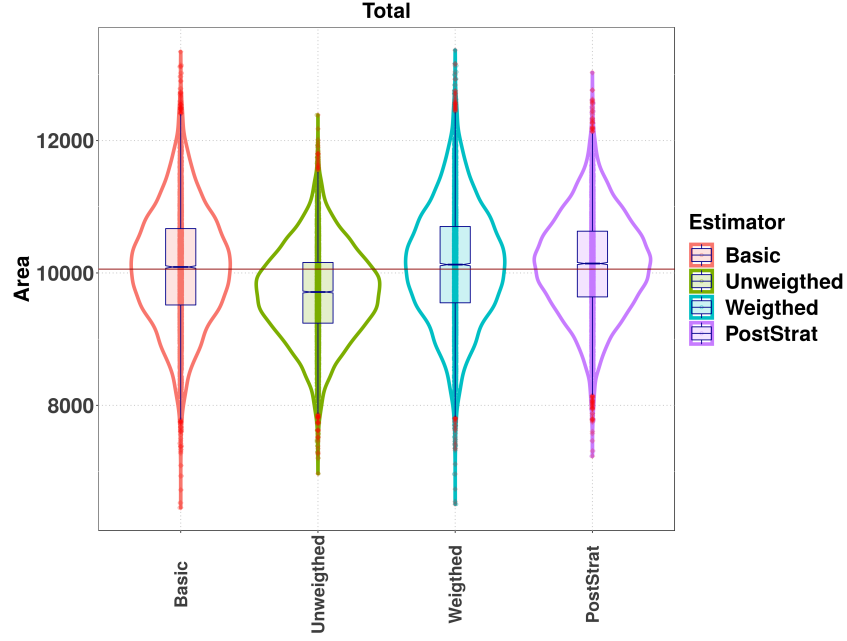


Figure 3: Violin-plots with Boxplots of Estimators Performances for 10000 Monte Carlo Replicate. Sample size: 120

Table 8: Summary Statistics of 10000 Monte Carlo Replications. Sample size: 180

Stratum	Estimator	Estimador	Bias	Rel.bias	EQM	SD
I	Basic	5021.72	28.07	0.01	234212.18	483.16
	Unweighed	6428.54	1434.89	0.29	2272620.21	462.32
	Weighthed	5981.36	987.71	0.20	1395449.61	648.01
	PostStrat	5983.84	990.19	0.20	1163961.53	428.38
II	Basic	3720.07	1.15	0.00	130199.04	360.85
	Unweighed	3005.05	-713.87	-0.19	609953.56	316.79
	Weighthed	3884.77	165.85	0.04	262296.38	484.58
	PostStrat	3886.67	167.75	0.05	196703.63	410.59
III	Basic	1351.38	5.92	0.00	109559.74	330.96
	Unweighed	275.86	-1069.60	-0.79	1145795.42	41.79
	Weighthed	236.61	-1108.85	-0.82	1231279.54	41.49
	PostStrat	236.51	-1108.95	-0.82	1231038.19	35.49
Total	Basic	10093.18	35.14	0.00	476483.75	689.42
	Unweighed	9709.45	-348.58	-0.03	420949.66	547.24
	Weighthed	10102.74	44.71	0.00	477090.42	689.30
	PostStrat	10107.02	48.98	0.00	344144.08	584.62

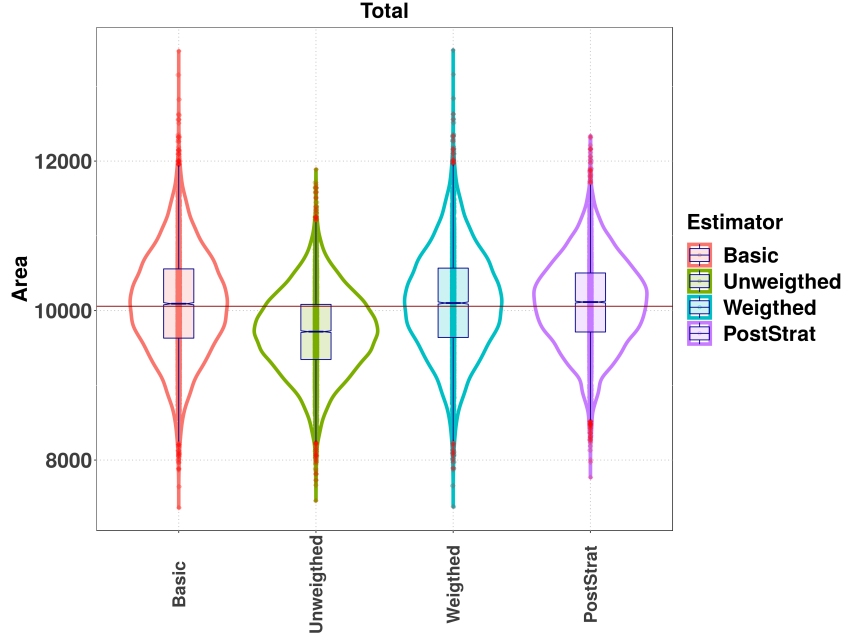


Figure 4: Violin-plots with Boxplots of Estimators Performances for 10000 Monte Carlo Replicate. Sample size: 180

## 6 Concluding remarks

The simulation investigated the effect of different estimators trying to make corrections to strata misclassification of segments in a problem of area sampling of square segments where stratification was made based on crowd-sourcing.

The post-stratified estimator shows better performance than the others cases, as expected.

However, no information regarding the actual size of each stratum is available in practice. In this situation, the investigation has shown that the performances of the basic and the weighted estimators are similar for producing estimates to the whole population. However, inside each stratum, the weighted estimators is subjected to bias that can be non negligible.

In addition, the unweighted estimators has shown a poor performance in all cases. Another factor that may affect estimators behavior is the percentage of misclassification. In all cases here, the percentages were kept the same as in the Goiana's experiment.

**Acknowledgements:** RO gratefully acknowledges from CNPq (Brazil).

## References

- Abreu, D. A., P. Arroway, A. C. Lamas, K. K. Lopiano, and L. J. Young (2010). Using the census of agriculture list frame to assess misclassification in the june area survey. In *Proceedings of the Joint Statistical Meetings*, pp. s/n.
- Ferraz, C., J. Delincé, and J. Gallego (2018). Agricultural Master Sampling Frames: Lessons learned from international field experiments and case studies. Technical Report No. 39GSARS Technical Report: Rome. Technical report, FAO.
- Jang, D., A. Sukasih, X. Lin, K. H. Kang, and S. H. Cohen (2009). Effects of misclassification of race/ethnicity categories in sampling stratification on survey estimates. In *Proceedings of the*

*American Statistical Association, Survey Methods Section*, pp. 3414–28. American Statistical Association Alexandria, VA.

Lamas, A. C., D. A. Abreu, P. Arroway, K. K. Lopiano, and L. J. Young (2010). Modeling misclassification in the june area survey. In *Proceedings of the section on survey research methods JSM*, pp. s/n.

Mulrow, J. and L. Woodburn (1990). *An Investigation of Stratification Errors*. American Statistical Association 1990 Proceedings of the Section on Survey Research Methods: ASA.

Table 9: Summary Statistics of 10000 Monte Carlo Replications. Sample size: 180

Stratum	Estimator	Mean	Bias	Relative bias	MSE	SD
I	Basic	5027.53	33.88	0.01	364673.14	602.96
	PostStrat	6462.59	1468.94	0.29	2487798.04	574.49
	Unweigthed	6796.53	1802.88	0.36	3616724.91	605.30
	Weigthed	6452.70	1459.05	0.29	4947341.77	1678.92
II	Basic	3728.02	9.10	0.00	460885.40	678.86
	PostStrat	4782.97	1064.05	0.29	1878180.97	863.74
	Unweigthed	3447.14	-271.78	-0.07	461695.72	622.79
	Weigthed	4716.95	998.03	0.27	3720555.89	1650.69
III	Basic	1333.91	-11.55	-0.01	3231841.80	1797.79
	PostStrat	403.83	-941.63	-0.70	895157.49	92.12
	Unweigthed	501.41	-844.05	-0.63	725494.98	114.35
	Weigthed	400.83	-944.63	-0.70	917662.48	159.16
Total	Basic	10089.46	31.43	0.00	3993771.56	1998.30
	PostStrat	11649.39	1591.36	0.16	3597847.92	1032.24
	Unweigthed	10745.09	687.05	0.07	1217802.85	863.62
	Weigthed	11570.48	1512.45	0.15	6098636.32	1952.31