



Sampling redesign of soil penetration resistance in spatial t-Student models

Leticia E. D. Canton¹, Luciana P. C. Guedes¹, Miguel A. Uribe-Opazo¹, Rosângela A. B. Assumpção² and Tamara C. Maltauro¹

¹ Western Paraná State University (UNIOESTE), 2069 Universitária Street, 85819-110, Cascavel, Paraná, Brazil. ² Federal Technological University of Paraná (UTFPR), 19 Cristo Rei Street, 85902-490, Toledo, Paraná, Brazil.

Abstract

Aim of study: To reduce the sample size in an agricultural area of 167.35 hectares, cultivated with soybean, to analyze the spatial dependence of soil penetration resistance (SPR) with outliers.

Area of study: Cascavel, Brazil

Material and methods: The reduction of sample size was made by the univariate effective sample size (ESS_t) methodology, assuming that the t-Student model represents the probability distribution of SPR.

Main results: The radius and the intensity of spatial dependence have an inverse relationship with the estimated value of the ESS_t . For the depths of SPR with spatial dependence, the highest estimated value of the ESS_t reduced the sample size by 40%. From the new sample size, the sampling redesign was performed. The accuracy indexes showed differences between the thematic maps with the original and reduced sampling designs. However, the lowest values of the standard error in the parameters of the spatial dependence structure evidenced that the new sampling design was appropriate. Besides, models of semivariance function were efficiently estimated, which allowed identifying the existence of spatial dependence in all depth of SPR.

Research highlights: The sample size was reduced by 40%, allowing for lesser financial investments with data collection and laboratory analysis of soil samples in the next mappings in the agricultural area. The spatial t-Student model was able to reduce the influence of outliers in the spatial dependence structure.

Additional key words: effective sample size; geostatistics; robust methods; simulation

Abbreviations used: CrVa (cross-validation); CV (coefficient of variation); ESS_n (univariate effective sample size of variables with normal probability distribution); ESS_t (univariate effective sample size of variables with Student's t-distribution); GPS (global positioning system); OA (overall accuracy); PA (precision agriculture); RNE (relative nugget effect); SD (standard deviation); SE (standard error); SPR (soil penetration resistance); T (Tau concordance index); UTM (Universal Transverse Mercator).

Authors' contributions: All authors: conceptualized the paper, statistical analysis of data, final revision and discussion. LEDC: reviewed the literature and edited the working versions of the manuscript.

Citation: Canton, LED; Guedes, LPC; Uribe-Opazo, MA; Assumpção, RAB; Maltauro, TC (2021). Sampling redesign of soil penetration resistance in spatial t-Student models. Spanish Journal of Agricultural Research, Volume 19, Issue 1, e0202. <https://doi.org/10.5424/sjar/2021191-16949>

Received: 20 May 2020. **Accepted:** 18 Mar 2021.

Copyright © 2021 INIA. This is an open access article distributed under the terms of the Creative Commons Attribution 4.0 International (CC-by 4.0) License.

Funding agencies/institutions

Coordination for the Improvement of Higher Education Personnel (CAPES). Financing Code 001

National Council for Scientific and Technological Development (CNPq)

Post-Graduate Program in Agricultural Engineering (PGEAGRI)

Competing interests: The authors have declared that no competing interests exist.

Correspondence should be addressed to Leticia E. D. Canton: leticiacanton@hotmail.com

Introduction

The Brazilian economy is directly related to agribusiness, and soybean (*Glycine max* (L.) Merrill) lead this scenario, which figures as the main grain exported by Brazil. Given the economic importance of this commodity, to

preserve the productivity and increase it, it is important to know the spatial variability of soybean yield and its relationship with the physical and chemical properties of the soil (Sobjak *et al.*, 2016). From this perspective, precision agriculture (PA) techniques use the knowledge of the spatial variability of grain yield and the physical and

chemical properties of the soil, to find the ideal application of the nutrient according to local needs (Molin *et al.*, 2015). The premise of PA is to use localized management of agricultural inputs to increase profits, reduce losses, and preserve the environment (Alamo *et al.*, 2012; Bier & Souza, 2017).

Geostatistics can help PA, as its techniques make it possible to determine the spatial dependence structure and describe the spatial variability of the yield of soybean and the soil attributes (Dalposso *et al.*, 2016, 2018; De Bastiani *et al.*, 2017; Schemmer *et al.*, 2017; Fagundes *et al.*, 2018; Grzegozewski *et al.*, 2020). The geostatistical techniques consider the value observed and geographic location of the physical-chemical properties of the soil, considering a sampling of some georeferenced points in the area. Thus, the entire area is characterized by a small representative portion of it (Wang *et al.*, 2013).

Knowing the spatial distribution of soil attributes and agricultural production is possible, even for small farmers. Combining sample planning and spatial statistics techniques, it is possible to characterize the spatial variability of attributes without using equipment with high investment, such as a harvest monitor (Schemberger *et al.*, 2017).

Also, to better understand the nutritional characteristics of the soil, it is important to combine samples of macro- and micro- nutrients and physical attributes, such as soil penetration resistance (SPR), which is related to the analysis of soil compaction. Compacted soils tend to hinder the availability of nutrients and water to the plant, which interferes with the growth of the roots and, consequently, with the development of the plant and the grain, thus affecting productivity (Valadão *et al.*, 2015, 2017; Marinello *et al.*, 2017; Sivarajan *et al.*, 2018; Colombi & Keller, 2019).

Still, in terms of sampling, there are studies that aim to reduce costs with collection and laboratory analysis of the sample. These studies proposed methods to reduce the number of sampling points to be used in future experiments in the agricultural area, without having a considerable loss in its mapping (Griffith, 2005; Guedes *et al.*, 2014, 2016; Domenech *et al.*, 2017; Maltauro *et al.*, 2019). One of the proposals is the effective sample size, which considers that some sample points may be highly correlated with each other, providing unnecessary cost with collection and laboratory analyzes, since such points add repeated information regarding spatial dependence (Vallejos & Osorio, 2014). The effective sample size represents the estimation of a new sample size considering the effects of the spatial autocorrelation and the purpose of estimating the sample mean of the value of the georeferenced variable as precisely as possible (Griffith, 2005).

The univariate effective sample size estimation developed by Griffith (2005), assumes that the georeferenced attribute has a normal probability distribution. However, there are georeferenced data that do not present a normal

probability distribution, especially because such distribution is sensitive to outliers (Fagundes *et al.*, 2018). In this way, Vallejos & Osorio (2014) suggested another more inclusive approach to calculate the estimated value of univariate effective sample size, which considers the presence of outliers and assumes that the georeferenced variable has Student's t-distribution. The Student's t-distribution allows the class of errors to be extended to other probability distributions to better accommodate the outliers (Assumpção *et al.*, 2014; De Bastiani *et al.*, 2015; Schemmer *et al.*, 2017).

The estimation of effective sample size requires an initial sampling design in the agricultural area and the knowledge of the spatial dependence structure of the georeferenced variable. Generally, when this information is not previously known and the data collection is being initiated, the initial sample size can be determined by the ratio of area to sample size (Wang *et al.*, 2013). For example, the PA recommend considering a maximum of two hectares per sampling point (Molin *et al.*, 2015).

Considering the availability of information obtained previously from the sample design in an experimental area, this study had as main objectives: i) to consider variables that present Student's t-distribution, using the expectation-maximization (EM) algorithm to model the data (Assumpção *et al.*, 2014); ii) to use the spatial dependence structure of SPR to redefine and to reduce the number of sample elements collected in this area by univariate effective sample methodology, considering the existence of the sample points correlated with each other.

Material and methods

We developed two studies: in the first, simulated data was considered, and in the second, we used data on SPR obtained in an agricultural area with soybean cultivation. The simulation study complements the agricultural one because with the simulated data is possible to reproduce a variety of scenarios present in the real data. Therefore, the two studies add practical and theoretical knowledge about sample resizing in soil attributes with spatial dependence structure.

Description of simulations

Consider a stochastic process $\{Y(\mathbf{s}_i), \mathbf{s}_i \in S \subset \mathbb{R}^2\}$, $i = 1, \dots, n$, stationary and isotropic, in which $\mathbf{Y} = (Y(\mathbf{s}_1), \dots, Y(\mathbf{s}_n))^T$ is a $n \times 1$ random vector, where $Y(\mathbf{s}_1), \dots, Y(\mathbf{s}_n)$ are the observed values of the random variable under study in \mathbf{s}_i sampled spatial locations, with $i = 1, \dots, n$ and $\mathbf{s}_i \in S \subset \mathbb{R}^2$. Suppose that \mathbf{Y} has an n -varied Student's t-distribution (De Bastiani *et al.*, 2015), *i.e.*, $\mathbf{Y} \sim t_n(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu)$, where $\boldsymbol{\mu}$ is the mean of Y , a constant

value in all \mathbf{s}_i spatial locations $\mathbf{1}$ is an $n \times 1$ unit-dimensional vector; ν ($\nu > 0$) is the degree of freedom fixed; and $\Sigma = \varphi_1 \mathbf{I}_n + \varphi_2 \mathbf{R}(\varphi_3)$ is an $n \times n$ covariance matrix, non-singular, where $\varphi_1 \geq 0$ and $\varphi_2 \geq 0$ are the nugget effect and partial sill parameters, respectively, \mathbf{I}_n is the $n \times n$ identity matrix, and $\mathbf{R}(\varphi_3)$ is an $n \times n$ symmetric matrix, where $\varphi_3 > 0$ is a function of the range ($g(\varphi_3) = a$). The practical range (a) is the spatial dependence radius, the distance at which spatial dependence exists between samples. The parameters φ_1 , φ_2 , and φ_3 , make up the spatial dependence structure of a georeferenced variable (Diggle & Ribeiro Jr, 2007; Soares, 2014). We considered 11 variables (V1, ..., V11) with different spatial dependence structures (Fig. 1A). The variables were obtained by simultaneously varying the spatial dependence radius (a) and the intensity of spatial dependence, measured by the relative nugget effect (RNE). As we set the parameter φ_2 value then the RNE was directly influenced by the variation of the nugget effect (φ_1). The smallest spatial dependence radius used was 0.3 km, and the largest ranged between 1.0 and 1.2 km. The remaining practical ranges (0.5 and 0.6 km) were considered intermediate based on the maximum distance from the agricultural area (1.8 km). The RNE was

considered from moderate (between 25% and 75%) to strong ($\leq 25\%$) (Cambardella *et al.*, 1994).

Given the linear spatial model (Uribe-Opazo *et al.*, 2012), we performed 100 simulations for each of the 11 variables using a Monte Carlo experiment from the Cholesky decomposition of the scale matrix Σ (Cressie, 2015). Each simulation generates a random sample set of these variables, maintaining the characteristics of the spatial dependence structure, and represents different datasets in different agricultural areas or crop years (Mooney, 1997). In these simulations, we fixed the degree of freedom ($\nu = 5$), the mean ($\mu = 5$), the partial sill ($\varphi_2 = 1$), and the exponential model. As sample planning for the simulations, we used the same configuration (lattice plus close pairs) from the commercial agricultural area under study (Fig. 1C). Other information about the simulations is given in the methodological scheme (Fig. 1B).

Following the scheme presented in Fig. 1B, after apply the EM algorithm to estimate the parameter vector θ for each simulated variable, the value of the effective sample size using the Student's t-distribution was estimated (ESS_t , Eq. 1) (Vallejos & Osorio, 2014):

$$\widehat{ESS}_t = \frac{\nu+n}{\nu+n+2} \mathbf{1}^T \mathbf{V}(\hat{\varphi})^{-1} \mathbf{1} \quad (1)$$

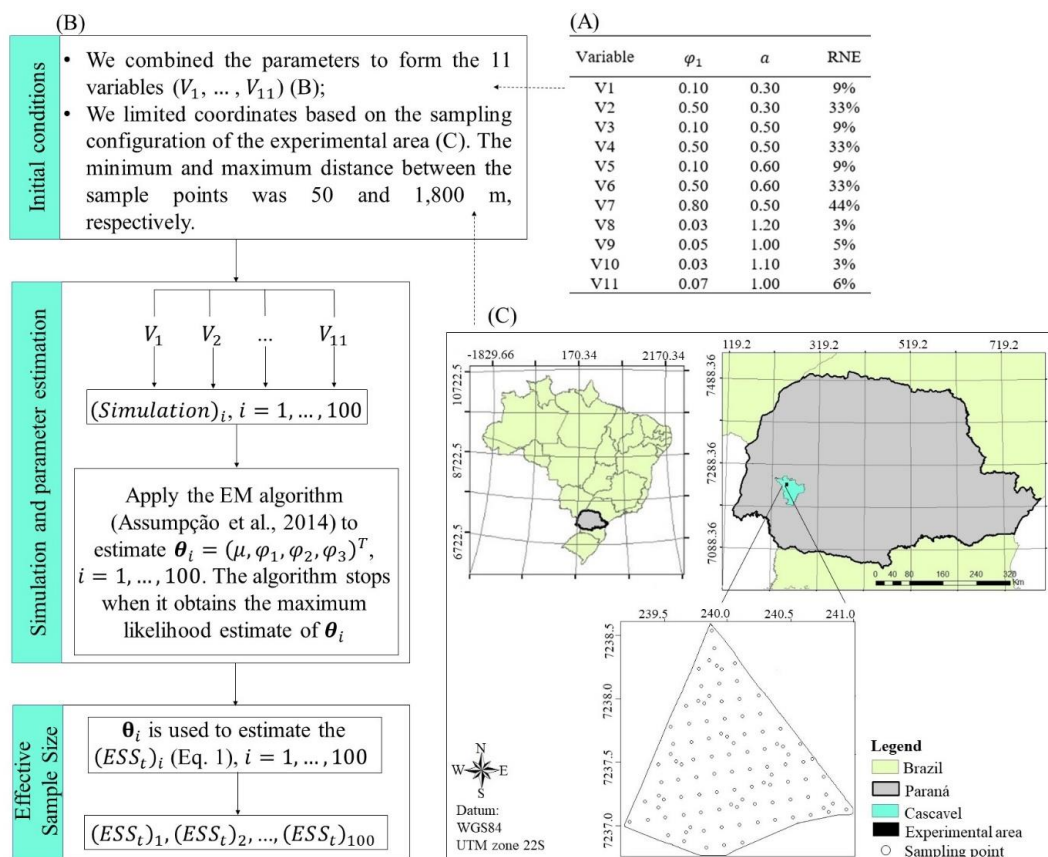


Figure 1. (A) Values of the parameters that define the spatial dependence structure of the simulated variables. (B) Methodological scheme used in simulation studies. (C) Experimental area with the location of the points sampled using UTM coordinates. μ : mean. $\hat{\varphi}_1$: nugget effect. $\hat{\varphi}_2$: partial sill. $\hat{\varphi}_3$: function of the range. a : practical range (kilometers). $RNE = 100 \frac{\hat{\varphi}_1}{\hat{\varphi}_1 + \hat{\varphi}_2}$: relative nugget effect (%). ESS_t : effective sample size.

$$v(\hat{\boldsymbol{\varphi}})_{ij} = \begin{cases} 1, & \text{se } i = j = 1, \dots, n \\ \frac{\widehat{\varphi}_2 r_{ij}}{\widehat{\varphi}_1 + \widehat{\varphi}_2}, & \text{se } i \neq j = 1, \dots, n \end{cases} \quad (2)$$

where n is the number of simulated sampling points in the original grid ($n \geq 1$); v is the degree of freedom ($v > 2$) $\mathbf{1}$ is an $n \times 1$ unit vector; $\mathbf{V}(\hat{\boldsymbol{\varphi}}) = [(v(\hat{\boldsymbol{\varphi}})_{ij})]$ is an $n \times n$ estimated spatial correlation matrix of the sample points, where the estimated spatial correlation between the i -th and the j -th sampling point are given by (Eq. 2); r_{ij} are the elements of the $\mathbf{R}(\varphi_3)$ matrix, which calculation depends on the geostatistical model and on the Euclidean distance between observations (De Bastiani *et al.*, 2015); and $\widehat{\varphi}_1$, $\widehat{\varphi}_2$ are the estimated values of the nugget effect and partial sill parameters, respectively.

What differs in estimating the univariate effective sample size by considering random vectors with normal probability distribution (ESS_n) in relation to those with Student's t-distribution (ESS_t), is the constant $\frac{(v+n)}{(v+n+2)}$. We obtained this constant from the Fisher information matrix for linear spatial models with Student's t-distribution (De Bastiani *et al.*, 2015). As $v > 2$ and $n \geq 1$, we have $v+n+2 > v+n$ and ESS_t is necessarily lower than ESS_n .

Description of the experimental data

The dataset comes from a commercial area with 167.35 hectares, cultivated with soybean, located in the municipality of Cascavel-Paraná-Brazil, with approximate geographical coordinates of latitude 24.95° South and longitude 53.37° West, and 650 m of average altitude (Fig. 1C). The climate of the region is temperate mesothermic and superhumid, climate type Cfa (Koeppen) (Aparecido *et al.*, 2016), with an average annual temperature of 21°C. The soil is classified as a Red Dystroferric Latosol with clay texture (EMBRAPA, 2013).

We used a lattice plus close pairs sampling design, with 102 sampling points. This design contained a regular grid (with minimum distance between points equals to 141 m), to which we added 19 sample points (locations). These added locations presented smaller distances with some points of the regular grid (50 m and 75 m). The sample was georeferenced and located with the aid of a signal receiving apparatus with a Geoexplore 3 (Trimble®) Global Positioning System (GPS) set up for the Universal Transverse Mercator (UTM) coordinate system.

In this study, soil resistance to root penetration (in MPa) at depths of 0-10 cm (SPR 0-10 cm), 11-20 cm (SPR 11-20 cm), 21-30 cm (SPR 21-30 cm), and 31-40 cm (SPR 31-40 cm) were used. In terms of improvement in soil management, the study of the spatial dependence of SPR has important agricultural relevance, since this soil attribute is inversely related to root growth and crop yield (Gül-

ser *et al.*, 2016). The experimental data of this physical attribute refers to the crop year 2015-2016 and belongs to the database of the Laboratory of Spatial Statistics and the Laboratory of Applied Statistics of the Western Paraná State University (UNIOESTE), Cascavel/Brazil.

The determination of SPR was measured by the penetrometer, as follows: for each sampling point, we performed three readings per centimeter, from 0 to 40 cm, covering the four depths considered (0-10 cm, 11-20 cm, 21-30 cm, and 31-40 cm). The data obtained was transformed in MPa, and the value of the SPR at each depth consisted of the arithmetic mean of the three measurements.

Soil penetration resistance was assumed to have a t-Student probability distribution. From the original sampling design and for each depth, we performed the exploratory and geostatistical analyzes of SPR (Figs. 2A and 2B, respectively). The analyses performed are described in the methodological scheme of Fig. 2, and more information about the methodology is obtained in Cressie (2015).

For each layer of SPR (at depths 0-10 cm, 11-20 cm, 21-30 cm, and 31-40 cm), the value of the effective sample size was estimated (ESS_t , Eq. 1) (Fig. 2) by the same methodology applied in the simulated data.

Through the estimated ESS_t values in each SPR layer, we redefined a single reduced sample size. The highest estimated value of the ESS_t was taken ($n^* = \text{MAX}(ESS_t)$, Fig. 2) from the variables with spatial dependence, *i.e.*, variables in which the value of the spatial dependence radius is not small compared relative to the size of the experimental area and which intensity of spatial dependence (RNE) was at least moderate (Cambardella *et al.*, 1994). We used only georeferenced variables with spatial dependence in the calculation of the ESS_t since georeferenced attributes without spatial dependence do not present a reduction in the number of sample points (Vallejos & Osorio, 2004).

The highest value criterion was established since a greater number of sampling points is better to capture the spatial variability of variables that have different spatial dependence structures (Pautsch *et al.*, 1998; Diggle & Ribeiro Jr, 2007). Therefore, the tendency is to obtain more representative thematic maps concerning the spatial variability of the attribute in experimental area (Kestring *et al.*, 2015). This is justified by two characteristics: (a) homogeneous variables (with less spatial variability in the area), can be collected with a smaller number of sample units, which would avoid redundant data or oversampling; and (b) variables with rapid change in spatial structure can be collected more intensively, which would avoid undersampling.

To verify the suitability of the reduced sample size, concerning the original sampling design (Fig. 1C), a random design of the original sampling design with sample size n^* was selected. For this reduced sample size, the exploratory and geostatistical analysis were also performed

(Figs. 2A and 2B, respectively). Finally, we compared the results obtained between the two sample configurations (original and reduced), using the methodologies presented in Fig. 2 (C and D).

The simulations, and the statistical and geostatistical analysis, were prepared in the software R (R Development Core Team, 2020) using the geoR package (Ribeiro Jr & Diggle, 2001). A computational routine developed in the software R (R Development Core Team, 2020) using the geoR (Ribeiro Jr & Diggle, 2001) and matrixcalc (Novomestky, 2012) packages (and available at goo.gl/JrvtnJ) to estimate the effective sample size (ESS_t).

Results

Simulation studies

The mean and the standard deviation of the estimated values of the ESS_t were similar for most pairs of variables in which the values of the nugget effect were different and the fixed range was maintained (V1 and V2; V3, V4, and V7; V5 and V6; V9 and V11) (Fig. 3). The estimated ESS_t

values evidenced the existence of three groups of variables (Fig. 3). The first two groups presented, respectively, the highest and intermediate estimated ESS_t values, being them: the group formed by variables V1 and V2, whose estimated mean value of the ESS_t was 40 and 44 sample points, in that order. Variables V3, V4, V5, V6, and V7 formed the second group, where the estimated mean value of the ESS_t ranged from 15 to 31 sample points. These two groups of variables also exhibited high values of standard deviations, with ESS_t variation of 11 to 14 sample points.

The simulated variables V1 and V2 have a small practical range ($\alpha = 0.3$ km), mainly when compared to the maximum distance between the coordinates of the simulated area (~ 1.8 km). Variables V3, V4, V5, V6, and V7 exhibited spatial dependence radius slightly higher than those of the first group (ranging from 0.5 to 0.6 km), which contributed to the fact that the estimated ESS_t values were smaller when compared to those obtained in the previous group. The third group, formed by variables V8, V9, V10, and V11, presented the smallest mean values of ESS_t (ranged from 6 to 8 sample points) (Fig. 3). These four variables have in common, the largest values of the simulated spatial dependence radius (between 1.0 and 1.2

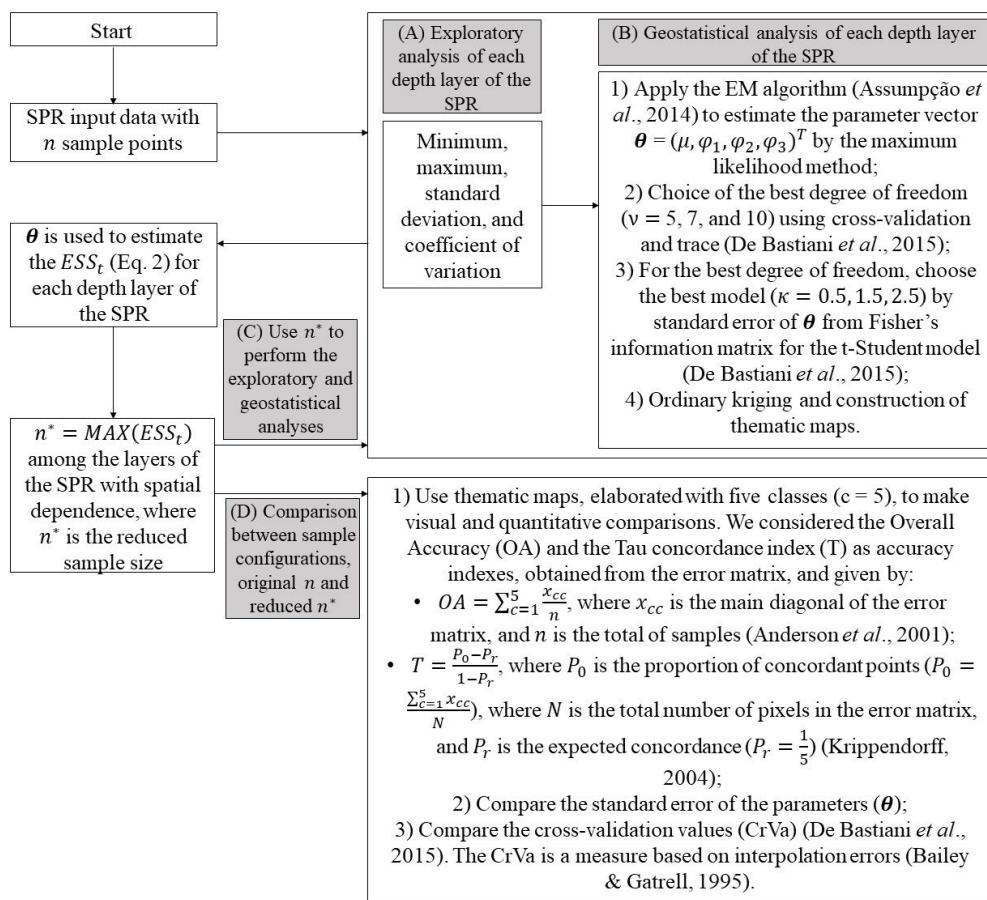


Figure 2. Methodological scheme used in experimental data. μ : mean. $\widehat{\varphi}_1$: nugget effect. $\widehat{\varphi}_2$: partial sill. $\widehat{\varphi}_3$: function of the range. ESS_t : effective sample size.

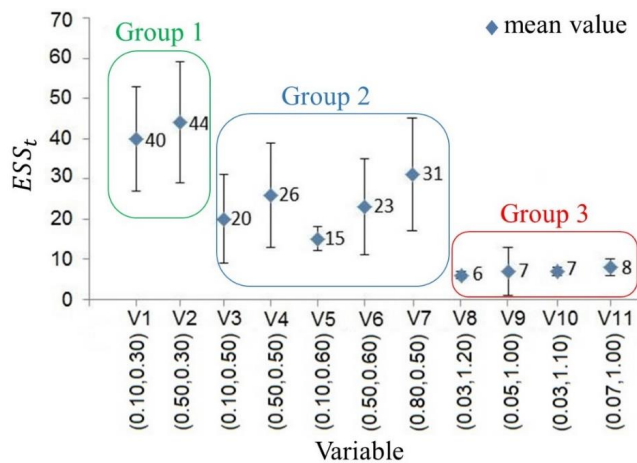


Figure 3. Mean and standard deviation of the estimated values of the effective sample size ESS_t for each variable, considering all the simulations. The parameters nugget effect ($\widehat{\varphi}_1$) and practical range (a), respectively, are shown in parentheses.

km). In general, the estimated ESS_t value ranged from 6 to 44 sample points and provided a reduction between 57% and 95% in the number of sampling points (Fig. 3).

Application of the methodology in soil penetration resistance

The estimated ESS_t value for SPR at depths 11-20 cm and 21-30 cm was 95 e 101, respectively. The SPR observed at depths 0-10 cm and 31-40 cm had higher reductions in the number of sampling points, with \widehat{ESS}_t equal to 51 and 60, respectively, which represents a reduction between 40% and 50%.

Considering the layers of SPR in which spatial dependence was identified (at depths 0-10 cm and 31-40 cm) and the maximum estimated value of the effective sample size observed in these layers, a sample resizing was obtained, reducing the sample size to 60 sample points. Thus, a new sample configuration with 60 points, chosen randomly from the 102 sample points of the original grid, was selected for the study of spatial dependence of SPR.

In the exploratory analysis, the values of the coefficient of variation (CV) showed that the SPR variability is greater

in the surface and it is reduced when increasing the sampling depth in the soil (Table 1). The magnitudes of the CVs indicated that there was a medium dispersion of SPR at all depths (Warrick & Nielsen, 1980) (Table 1). Besides, we observed that the reduction in the number of sample points did not influence the SPR variability (Table 1).

The depths 11-20 cm (Fig. 4B) and 31-40 cm (Fig. 4D) showed the greatest amount of outliers (four each), located in the central and western regions of experimental area. The sample points 82 and 34 exhibited outliers, in all depth layers of SPR, except at depth 31-40 cm, in which point 34 was not considered an outlier.

We observed in the geostatistical analysis that, for all depth layers of SPR, the spatial dependence structure can be considered isotropic, *i.e.*, depends only on the distance separating the locations observed, and does not differ with the direction (Guedes *et al.*, 2013).

The results about the best values for the degree of freedom (ν) and the shape parameter κ (Table 2), showed that for both sample sizes, the model and degree of freedom were the same only for SPR at depth of 31-40 cm. We verified in this depth the lowest values of the standard error (SE) in the Matérn family model with $\kappa = 0.5$, for the degree of freedom $\nu = 10$. For SPR at depth of 0-10 cm, the lowest values estimated from the SE were found in the Matérn family model with $\kappa = 2.5$ and shape parameter $\nu = 5$, or the original sampling design; and with $\kappa = 0.5$ and $\nu = 10$ for the reduced sampling design. At depth 11-20 cm, in both sampling designs, was adjusted the Matérn family model $\kappa = 2.5$ to the semivariance function, but with different degrees of freedom $\nu = 5$ for the original grid, and $\nu = 10$ for the reduced grid.

Finally, at depth 21-30 cm of the SPR, although the sampling designs presented the same value for the shape parameter ($\nu = 5$), the lowest estimated values of the SE were obtained by the Matérn family model with $\kappa = 1.5$ and 2.5, respectively, for the original and reduced sampling designs (Table 2). The estimated values of these SEs, at depths 0-10 cm and 11-20 cm of the SPR (Table 2), were smaller in the estimated models considering the reduced sampling design when compared to values obtained in the estimated models using the original sampling design. Besides, for the other depth layers, the estimated value of

Table 1. Descriptive statistics of the four depth layers of soil penetration resistance (SPR, in MPa), considering the original ($n=102$ points) and reduced ($n^*=60$ points) sampling designs.

Attribute	Minimum		Mean		Maximum		SD		CV	
	n	n^*	n	n^*	n	n^*	n	n^*	n	n^*
SPR 0-10 cm	0.75	1.49	3.12	3.18	6.67	6.67	0.97	1.04	31.22	32.69
SPR 11-20 cm	1.78	1.78	3.01	3.05	5.53	5.53	0.62	0.67	20.70	22.12
SPR 21-30 cm	1.50	1.50	2.03	2.04	3.86	3.86	0.33	0.37	16.27	18.20
SPR 31-40 cm	1.48	1.56	2.08	2.00	3.95	3.95	0.37	0.42	18.01	19.83

SD: standard deviation. $CV = 100 \frac{SD}{Mean}$: coefficient of variation (%).

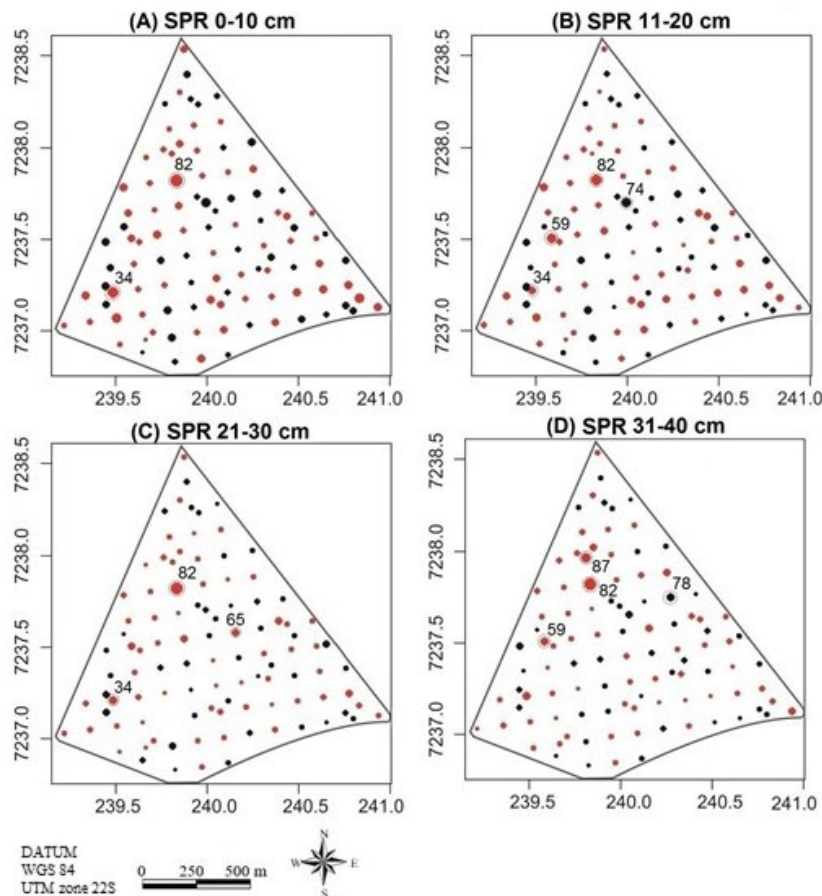


Figure 4. Post-plot graphic for all depth layers of soil penetration resistance (SPR). The black dots represent the 102 observations in the experimental area. The red dots indicate the spatial location of the 60 selected points to compose the reduced sampling design. The dot size is proportional to the measured value of the SPR at each sampling point. The numerate points indicate spatially the outliers in the experimental area.

the SE of the range function ($\widehat{\varphi}_3$), was also lower for the estimated models considering the reduced sample configuration (Table 2).

The values obtained by the cross-validation method showed a small increase in errors of the spatial prediction with the reduced sampling design (Table 2). The errors increased by 7.5%, 6.2%, 9.6%, and 11.5%, respectively, at depths 0-10 cm, 11-20 cm, 21-30 cm, and 31-40 cm of the SPR, comparing with the original sampling design.

For the original sampling design, the spatial dependence structure in the intermediate depth layers of the SPR (11-20 cm and 21-30 cm) presented pure nugget effect, due to the low values of the practical range (180.5 and 110.2 m) and the low spatial dependence ($RNE \geq 5\%$; Cambardella *et al.*, 1994) (Table 2).

Considering the reduced sampling design, the intensity of spatial dependence was moderate in these intermediate-depth layers (RNE between 25% and 75%; Cambardella *et al.*, 1994). Also, at depth 11-20 cm, there was an increase in the estimated value of the spatial dependence radius to the original sampling design (from 180.5 to 209.1 m). In the other depth layers of the SPR, there was a decrease

in the estimated practical range (ranging from 7.6 to 31.9 m), compared to that obtained with the original sampling design (Table 2).

The estimated values of the practical range were relatively low for both sampling designs and in all depth layers of the SPR. That is because the maximum distance in the experimental area is approximately 1,800 m, the ranges ranged from 110.2 to 291.5 m in the original sampling design, and from 78.3 to 273.2 m in the reduced sampling design (Table 2).

We observed visual differences between the maps elaborated considering the two sampling designs, which are most noticeable at depth 11-20 cm (Fig. 5B). According to the classification of Anderson *et al.* (2001), in most of the depth layers of the SPR, there was a low percentage of hits between the reference map (original sampling design) and the model map (reduced sampling design), because the estimated value of the overall accuracy (OA) was lower than 85%. This indicates that a smaller number of pixels were classified in the same class interval in both maps, evidencing differences between the elaborated maps considering the two sampling designs. The only

Table 2. Estimated values of the parameters that define the spatial dependence structure in each depth layer of soil penetration resistance (SPR, in MPa) from the best values of the shape parameters \mathcal{K} and ν considering the original (n = 102 points) and reduced (n* = 60 points) samplings designs.

		Estimated parameter								
		ν	\mathcal{K}	$\hat{\mu}$	$\hat{\varphi}_1$	$\hat{\varphi}_2$	$\hat{\varphi}_3$	\hat{a}	\overline{RNE}	CrVa
SPR 0-10 cm	n	5	2.5	3.1188 (0.1363)	0.6246 (0.2339)	0.3072 (0.3667)	0.0492 (0.0107)	0.2915	67.03%	0.8791
	n*	10	0.5	3.1751 (0.2412)	0.3447 (0.1275)	1.4263 (0.2447)	0.0911 (0.0027)	0.2731	19.46%	1.0219
SPR 11-20 cm	n	5	2.5	2.9992 (0.0571)	0.2842 (0.1087)	0.0209 (0.9389)	0.0305 (0.6621)	0.1805	93.14%	0.3922
	n*	10	2.5	3.0447 (0.1090)	0.2356 (0.0672)	0.2559 (0.2468)	0.0353 (0.0051)	0.2091	47.94%	0.4437
SPR 21-30 cm	n	5	1.5	2.0180 (0.0298)	0.0820 (0.0426)	0.0056 (3.3147)	0.0232 (6.4973)	0.1102	93.60%	0.1192
	n*	5	2.5	2.0382 (0.0431)	0.0707 (0.3244)	0.0356 (4.5560)	0.0132 (0.6775)	0.0783	66.50%	0.1446
SPR 31-40 cm	n	10	0.5	2.0659 (0.0485)	0.0635 (0.0166)	0.0755 (0.2478)	0.0742 (0.0483)	0.2225	45.67%	0.1350
	n*	10	0.5	2.0974 (0.0807)	0.1096 (0.0339)	0.1597 (0.2609)	0.0716 (0.0314)	0.2149	40.71%	0.1700

ν : degree of freedom. \mathcal{K} : shape parameter of the Matérn family model. Estimated values of: $\hat{\mu}$: mean, $\hat{\varphi}_1$: nugget effect, $\hat{\varphi}_2$: partial sill, $\hat{\varphi}_3$: function of the range, \hat{a} : practical range (kilometers), $\overline{RNE} = 100 \frac{\hat{\varphi}_1}{\hat{\varphi}_1 + \hat{\varphi}_2}$: relative nugget effect (%), CrVa: cross-validation. The estimated values of the standard error (SE) for each parameter are shown in parentheses.

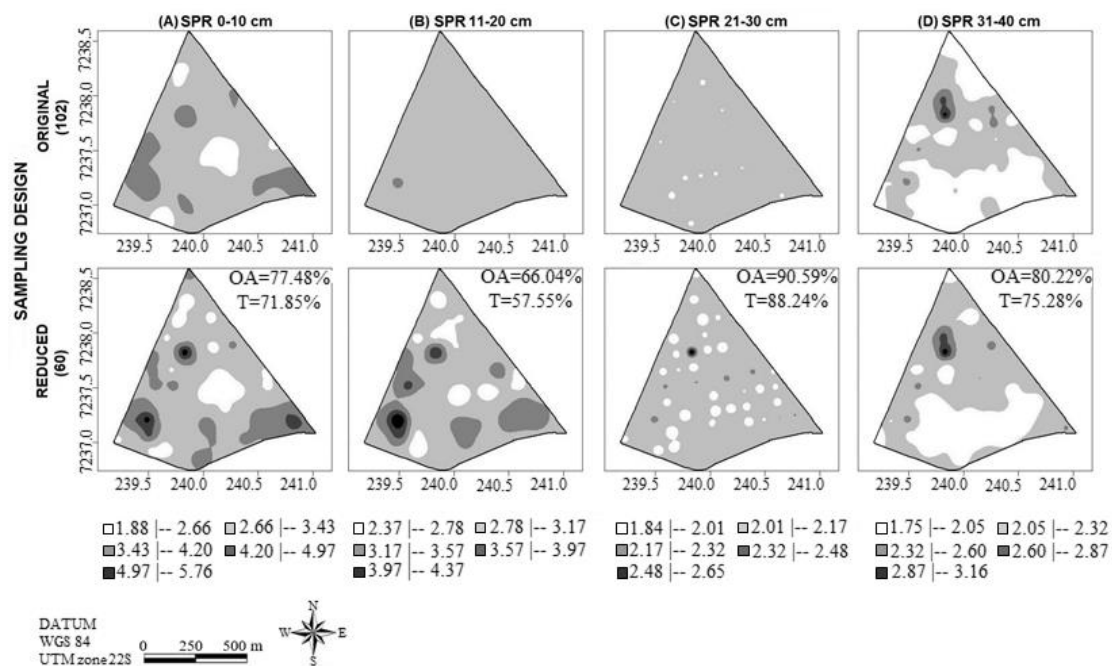


Figure 5. Thematic map of the estimated values for all depth layers of soil penetration resistance (SPR), in MPa, considering the original and reduced sampling designs with the same class intervals. Estimated values of OA (overall accuracy) and T (Tau concordance index) (%).

exception was at depth 21-30 cm (Fig. 5C) in which similarity between the maps made with the original and reduced sampling design was observed ($OA > 85\%$).

The Tau concordance index (T), unlike the OA, accounts for not only the proportion of pixels classified in the same class interval in the reference and model maps but also those whose classification was not the same in both maps. The maps made considering the original and reduced sampling designs presented from low to medium accuracy ($T < 0.80$; Krippendorff, 2004), with the exception at depth 21-30 cm of the SPR (Fig. 5).

Some classes of thematic maps elaborated using the original sampling design presented some null pixels, as can be seen visually at depths of 0-10 cm, 11-20 cm, and 21-30 cm of the SPR (Figs. 5A, 5B, and 5C, respectively). Besides, at depth 21-30 cm, where was obtained high values of the EG and Tau accuracy indexes, a high number of pixels (more than 90% of total) in the same classes was observed (Fig. 5C). Still about at depth 21-30 cm, we observed the formation of circular regions around the sample points (Fig. 5C).

Finally, the estimated values of SPR, using the reduced sampling design, showed the existence of limitations to root growth that varied from low to moderate in almost all the agricultural areas (Canarache, 1991).

Discussion

Simulation studies

Considering the 100 simulations of each variable, the graph with the means and standard deviations of the estimated values of the univariate effective sample size showed that the variation of the value of the nugget effect did not generate a relevant change in the estimated value of the effective sample size (Fig. 3). The practical range negatively influenced the estimated ESS_t values, since the greater the practical range of the variable, the lower the estimated ESS_t value. Although a different sample configuration and size, and even another probabilistic distribution (normal) were considered, the simulation studies of Vallejos & Osorio (2014) and Dal Canton *et al.* (2021) reached similar conclusions regarding the influence of the practical range in reducing the number of sampling points.

The high difference in the estimated ESS_t values (Fig. 3) can be explained by the discrepancy between the variables concerning the values of the parameters of spatial dependence, mainly regarding the practical range, which variation was of 0.3 to 1.2 km.

Studies carried out in agricultural areas of a smaller size than the one considered in this paper (< 50 ha), characterized the spatial dependence on soil attributes using a sample size smaller than 50 sample points (Carvalho *et al.*, 2013; Araújo *et al.*, 2014; Tavares *et al.*, 2014), sam-

ple size similar to that obtained in the present study for most simulated variables.

Application of the methodology in soil penetration resistance

The values of the CV obtained by Johann *et al.* (2004) and Bazzi *et al.* (2013) showed results similar to those of this work, with a moderate classification for CVs in agricultural areas in Western Paraná with soybean planting, and with similar conditions of management, climate, and soil. Besides, the SPR variability is reduced when increasing the sampling depth in the soil, corroborating with that obtained in this study.

The SPR at depths 11-20 cm and 21-30 cm practically did not show a reduction in sample size. This fact is justifiable, mainly due to the practical range influence on the estimated ESS_t value, verified in the simulation studies of the present study, and verified in Vallejos & Osorio (2014) and Dal Canton *et al.* (2021) as well. The SPR at depths 11-20 cm and 21-30 cm presented a small estimated value of the spatial dependence radius (110.2 and 180.5 m), relative to the size of the experimental area, and also a low intensity of spatial dependence ($\widehat{RNE} > 75\%$; Cambardella *et al.*, 1994) (Table 2). The higher reductions in the number of sampling points presented at depths 0-10 cm and 31-40 cm are due to the higher estimated values of practical range (291.5 and 222.5 m) (Table 2).

Griffith (2005) obtained a reduction in sample size (from 36% to 45%) similar to that found in this study, which varied between 40% and 50%, using different sample configurations, attribute probability distribution, and soil chemical attributes. Domenech *et al.* (2017) considered auxiliary information measurement to map the attribute of interest (soil depth to the petrocalcic horizon), and obtained a reduction in sample size similar to this study also (from 50% to 70%), although their methodologies for optimization and selection of sampling points were different from this study. The mentioned authors obtained sample reductions similar to those of the present study, and the thematic maps obtained by them were considered efficient.

It was found in the literature researches to analyze the spatial variability at different depths of SPR and used between 49 and 60 sample points (Rosalen *et al.*, 2011; Rodrigues *et al.*, 2014; Tavares *et al.*, 2014). Considering the new sample configuration, reduced to 60 sample points, these authors used values of sample size similar to this research, although the magnitude of their mapped experimental areas was lower than that of this study (< 50 ha).

Comparing the original and reduced sampling designs, the estimated values of the relative nugget effect (RNE) and the practical range indicate that with a reduced

number of sample points there was an increase in spatial dependence and minor changes in the spatial dependence radius (Table 2). Besides, the estimated values of the SEs of the estimated parameters that define the spatial dependence structure were smaller in the estimated models from the reduced sampling design for the majority of the cases (Table 2). This shows that even with a smaller number of sample points, it was possible to verify the existence of spatial dependence in all depth layers of SPR.

The increase in spatial prediction errors after sampling redesign was already expected, as the number of sample points was reduced by 40%. The literature shows that the greater the number of sample points, the better the result of the interpolation, as shown by the studies by Coelho *et al.* (2009), Kestring *et al.* (2015), and Guedes *et al.* (2016), using different sample densities and metrics to calculate errors. However, the greater the number of observations, the greater the financial cost. Thus, with the magnitude of sample reduction obtained in this study, the increase in spatial prediction errors can be considered small.

The results also indicate that even using a smaller number of sampling points in the study area, efficiently models were estimated to the semivariance function and that they were able to identify the existence of spatial dependence in all depth layers of SPR. This is an important feature of this study since the reduction of the sample size difficult the semivariance calculation (Kestring *et al.*, 2015).

However, although it is possible to verify the existence of spatial dependence in all depth layers of SPR, the visual analysis and the accuracy indices (OA and Tau) showed that there are differences between the thematic maps generated with the original and reduced sampling design, indicating that there was the influence of sample size on the spatial dependence characterization of SPR.

About the circular regions around the sampling points, identified at depth 21-30 cm of the SPR map (Fig. 5C), we observed the low estimated value of the practical range ($\hat{\alpha} = 78.30$ m), near the shortest distance between sample points (~ 50 m), which resulted in the formation of small subregions centered in the sample points, a phenomenon known as ‘bull eyes effect’ (Menezes *et al.*, 2016), and also observed by Dalposso *et al.* (2018) and Dal Canton *et al.* (2021).

The results showed that the univariate ESS_t methodology proved to be advantageous considering the lowest cost in the sampling process due to the 40% reduction in the sample size and the results obtained in the characterization of the spatial dependence in the experimental area. Also, the method proposed in this study obtained a single sample size for all attributes, based on the variables with spatial dependence structure and the maximum estimated value of the ESS_t among them.

References

- Alamo S, Ramos MI, Feito FR, Cañas JA, 2012. Precision techniques for improving the management of the olive groves of southern Spain. *Span J Agric Res* 10 (3): 583-595. <https://doi.org/10.5424/sjar/2012103-361-11>
- Anderson JR, Hardy EE, Roach JT, Witmer RE, 2001. A land use and land cover classification system for use with remote sensor data. U.S. Government Print Office, Washington DC. 41 pp.
- Aparecido LEO, Rolim GS, Richetti J, Souza PS, Johann JA, 2016. Köppen, Thornthwaite and Camargo climate classifications for climatic zoning in the State of Paraná, Brazil. *Cienc Agrotec* 40 (4): 405-417. <https://doi.org/10.1590/1413-70542016404003916>
- Araújo DR, Mion RL, Sombra WA, Andrade RR, Amorim MQ, 2014. Variabilidade espacial de atributos físicos em solo submetido à diferentes tipos de uso e manejo. *Rev Caatinga* 27: 101-115.
- Assumpção RAB, Uribe-Opazo MA, Galea M, 2014. Analysis of local influence in geostatistics using Student's t-distribution. *J Appl Stat* 41: 2323-2341. <https://doi.org/10.1080/02664763.2014.909793>
- Bailey TC, Gatrell AC, 1995. Interactive spatial data analysis. Longman Scientific & Technical, Essex. 432 pp.
- Bazzi CL, Souza EG, Uribe-Opazo MA, Nóbrega LH, Rocha DM, 2013. Management zones definition using soil chemical and physical attributes in a soybean area. *Eng Agríc* 33 (5): 952-964. <https://doi.org/10.1590/S0100-69162013000500007>
- Bier AB, Souza EG, 2017. Interpolation selection index for delineation of thematic maps. *Comput Electron Agric* 136: 202-209. <https://doi.org/10.1016/j.compag.2017.03.008>
- Cambardella CA, Moorman T, Parkin T, Karlen D, Novak J, Turco R, Konopka A, 1994. Field-scale variability of soil properties in central Iowa soils. *Soil Sci Soc Am J* 58: 1501-1511. <https://doi.org/10.2136/sssaj1994.03615995005800050033x>
- Canarache A, 1991. Factors and indices regarding excessive compactness of agricultural soils. *Soil Till Res* 19: 145-164. [https://doi.org/10.1016/0167-1987\(91\)90083-A](https://doi.org/10.1016/0167-1987(91)90083-A)
- Carvalho LCC, Silva FM, Araújo G, Ferraz S, Silva FC, Stracieri J, 2013. Variabilidade espacial de atributos físicos do solo e características agrônômicas da cultura do café. *Coffee Sci* 8: 265-275.
- Coelho EC, Souza EGD, Uribe-Opazo MA, Pinheiro Neto R, 2009. Influência da densidade amostral e do tipo de interpolador na elaboração de mapas temáticos. *Acta Sci Agron* 31 (1): 165-174. <https://doi.org/10.4025/actas-ciagron.v31i1.6645>
- Colombi T, Keller T, 2019. Developing strategies to recover crop productivity after soil compaction - A plant eco-physiological perspective. *Soil Till Res* 191: 156-161. <https://doi.org/10.1016/j.still.2019.04.008>

- Cressie NAC, 2015. *Statistics for spatial data*, rev. ed. John Wiley & Sons, NY. 928 pp.
- Dal Canton LE, Guedes LPC, Uribe-Opazo MA, 2021. Reduction of sample size in the soil physical-chemical attributes using the multivariate Effective Sample Size. *J Agric Stud* 9 (1): 357-376. <https://doi.org/10.5296/jas.v9i1.17473>
- Dalposso GH, Uribe-Opazo MA, Johann JA, 2016. Soybean yield modeling using bootstrap methods for small samples. *Span J Agric Res* 14 (3): e0207. <https://doi.org/10.5424/sjar/2016143-8635>
- Dalposso GH, Uribe-Opazo MA, Johann JA, Galea M, De Bastiani F, 2018. Gaussian spatial linear model of soybean yield using bootstrap methods. *Eng Agríc* 38 (1): 110-116. <https://doi.org/10.1590/1809-4430-eng.agric.v38n1p110-116/2018>
- De Bastiani F, Cysneiros AFJ, Cysneiros AHM, Uribe-Opazo MA, Galea M, 2015. Influence diagnostics in elliptical spatial linear models. *Test* 24: 322-340. <https://doi.org/10.1007/s11749-014-0409-z>
- De Bastiani F, Galea M, Cysneiros AHMA, Uribe-Opazo MA, 2017. Gaussian spatial linear models with repetitions: An application to soybean productivity. *Spat Stat* 21: 319-335. <https://doi.org/10.1016/j.spasta.2017.07.013>
- Diggle P, Ribeiro Jr PJ, 2007. *Model-based geostatistics*. Springer, Lancaster. 228 pp. <https://doi.org/10.1007/978-0-387-48536-2>
- Domenech MB, Castro-Franco M, Costa JL, Amiotti NM, 2017. Sampling scheme optimization to map soil depth to petrocalcic horizon at field scale. *Geoderma* 290: 75-82. <https://doi.org/10.1016/j.geoderma.2016.12.012>
- EMBRAPA, 2013. *Sistema brasileiro de classificação de solos*, 3ed. Empresa Brasileira de Pesquisa Agropecuária, Centro Nacional de Pesquisa de Solos, Brasília. 306 pp.
- Fagundes RS, Uribe-Opazo MA, Guedes LPC, Galea M, 2018. Slash spatial linear modeling: soybean yield variability as a function of soil chemical properties. *Rev Bras Cienc Solo* 42: 1-14. <https://doi.org/10.1590/18069657rbc20170030>
- Griffith DA, 2005. Effective geographic sample size in the presence of spatial autocorrelation. *Ann Am Assoc Geogr* 95: 740-760. <https://doi.org/10.1111/j.1467-8306.2005.00484.x>
- Grzegozewski DM, Cima EG, Uribe-Opazo MA, Guedes LPC, Johann JA, 2020. Spatial and multivariate analysis of soybean yield in the state of Paraná-Brazil. *J Agric Stud* 8 (1): 387-412. <https://doi.org/10.5296/jas.v8i1.16303>
- Guedes LPC, Uribe-Opazo MA, Ribeiro Jr PJ, 2013. Influence of incorporating geometric anisotropy on the construction of thematic maps of simulated data and chemical attributes of soil. *Chil J Agric Res* 73 (4): 414-423. <https://doi.org/10.4067/S0718-58392013000400013>
- Guedes LPC, Uribe-Opazo MA, Ribeiro Jr PJ, 2014. Optimization of sample design sizes and shapes for regionalized variables using simulated annealing. *Cienc Invest Agrar* 41 (1): 33-48. <https://doi.org/10.4067/S0718-16202014000100004>
- Guedes LPC, Ribeiro Jr PJ, Uribe-Opazo MA, De Bastiani F, 2016. Soybean yield maps using regular and optimized sample with different configurations by simulated annealing. *Eng Agríc* 36 (1): 114-125. <https://doi.org/10.1590/1809-4430-Eng.Agric.v36n1p114-125/2016>
- Gülser C, Ekberli I, Candemir F, Demir Z, 2016. Spatial variability of soil physical properties in a cultivated field. *Euras J Soil Sci* 5 (3): 192-200. <https://doi.org/10.18393/ejss.2016.3.192-200>
- Johann JA, Uribe-Opazo MA, Souza EGD, Rocha JV, 2004. Variabilidade espacial dos atributos físicos do solo e da produtividade em um Latossolo Bruno distrófico da região de Cascavel, PR. *Rev Bras Eng Agríc Ambient* 8 (2-3): 212-219. <https://doi.org/10.1590/S1415-43662004000200008>
- Kestring F, Guedes LPC, De Bastiani F, Uribe-Opazo MA, 2015. Thematic maps comparison of different sampling grids for soybean productivity. *Eng Agríc* 35 (4): 733-743. <https://doi.org/10.1590/1809-4430-Eng.Agric.v35n4p733-743/2015>
- Krippendorff K, 2004. *Content analysis: an introduction to its methodology*. Sage Publications, Beverly Hills. 412 pp.
- Maltauro TC, Guedes LPC, Uribe-Opazo MA, 2019. Reduction of sample size in the analysis of spatial variability of nonstationary soil chemical attributes. *Eng Agríc* 39: 56-65. <https://doi.org/10.1590/1809-4430-eng.agric.v39nep56-65/2019>
- Marinello F, Pezzuolo A, Cillis D, Chiumenti A, Sartori L, 2017. Traffic effects on soil compaction and sugar beet (*Beta vulgaris* L.) taproot quality parameters. *Span J Agric Res* 15 (1): e0201. <https://doi.org/10.5424/sjar/2017151-8935>
- Menezes MD, Silva SHG, Mello CR, Owens PR, Curi N, 2016. Spatial prediction of soil properties in two contrasting physiographic regions in Brazil. *Sci Agric* 73 (3): 274-285. <https://doi.org/10.1590/0103-9016-2015-0071>
- Molin JP, Amaral LR, Colaço AF, 2015. *Agricultura de precisão*. Oficina de Textos, São Paulo. 224 pp.
- Mooney CZ, 1997. *Monte Carlo simulation*. Sage Publications, Thousand Oaks. 112 pp. <https://doi.org/10.4135/9781412985116>
- Novomestky F, 2012. *matrixcalc: collection of functions for matrix calculations*. R package version 3.3.1. <https://cran.r-project.org/web/packages/matrixcalc/index.html>
- Pautsch GR, Babcock BA, Breidt FJ, 1998. *Optimal sampling under a geostatistical model*. Center for Agricultural and Rural Development, Iowa, USA. 32 pp.

- R Development Core Team, 2020. R: A language and environment for statistical computing. Version 4.0.0. R Foundation for Statistical Computing, Vienna, Austria.
- Ribeiro Jr PJ, Diggle PJ, 2001. geoR: a package for geostatistical analysis. *R News* 1: 15-18. <https://cran.r-project.org/web/packages/geoR/index.html>
- Rodrigues MS, Ramos RRD, Azevedo TP, Patrocínio Filho AP, Oliveira LG, 2014. Variabilidade espacial da resistência do solo à penetração em área capineira irrigada no semiárido. *Agropecuária Científica no Semiárido* 10: 161-166.
- Rosalen DL, Rodrigues MS, Chioderoli CA, Brandão FJC, Siqueira DS, 2011. GPS receivers for georeferencing of spatial variability of soil attributes. *Eng Agríc* 31 (6): 1162-1169. <https://doi.org/10.1590/S0100-69162011000600013>
- Schemberger EE, Fontana FS, Johann JA, Souza EG, 2017. Data mining for the assessment of management areas in precision agriculture. *Eng Agríc* 37 (1): 185-193. <https://doi.org/10.1590/1809-4430-eng.agric.v37n1p185-193/2017>
- Schemmer RC, Uribe-Opazo MA, Galea M, Assumpção RAB, 2017. Spatial variability of soybean yield through a reparametrized t-Student model. *Eng Agríc* 37 (4): 760-770. <https://doi.org/10.1590/1809-4430-eng.agric.v37n4p760-770/2017>
- Sivarajan S, Maharlooei M, Bajwa SG, Nowatzki J, 2018. Impact of soil compaction due to wheel traffic on corn and soybean growth, development and yield. *Soil Till Res* 175: 234-243. <https://doi.org/10.1016/j.still.2017.09.001>
- Soares A, 2014. *Geoestatística para ciências da terra e do ambiente*, 3rd ed. Press, Lisboa. 214 pp.
- Sobjak R, Souza EG, Bazzi CL, Uribe-Opazo MA, Betzek NM, 2016. Redundant variables and the quality of management zones. *Eng Agríc* 36 (1): 78-93. <https://doi.org/10.1590/1809-4430-Eng.Agric.v36n1p78-93/2016>
- Tavares UE, Montenegro AAA, Rolim MM, Silva JS, Vicente TFS, Andrade CWL, 2014. Variabilidade espacial da resistência à penetração e da umidade do solo em Neossolo Flúvico. *Water Resour Irrig Manage* 3 (2): 79-89. <https://doi.org/10.19149/2316-6886/wrim.v3n2p79-89>
- Uribe-Opazo MA, Borssoi JA, Galea M, 2012. Influence diagnostics in Gaussian spatial linear models. *J Appl Stat* 39: 615-630. <https://doi.org/10.1080/02664763.2011.607802>
- Valadão FCA, Weber OLS, Júnior DDV, Scapinelli A, Deina FR, Bianchini A, 2015. Adubação fosfatada e compactação do solo: sistema radicular da soja e do milho e atributos físicos do solo. *Rev Bras Cienc Solo* 39 (1): 243-255. <https://doi.org/10.1590/01000683rbcs20150144>
- Valadão FCDA, Weber OLS, Júnior DDV, Santin MFM, Scapinelli A, 2017. Teor de macronutrientes e produtividade da soja influenciados pela compactação do solo e adubação fosfatada. *Rev Ciênc Agrár* 40 (1): 183-195. <https://doi.org/10.19084/RCA15092>
- Vallejos R, Osorio F, 2014. Effective sample size of spatial process models. *Spat Stat* 9: 66-92. <https://doi.org/10.1016/j.spasta.2014.03.003>
- Wang JF, Jiang CS, Hu MG, Cao ZD, Guo YS, Li LF, Liu TJ, Meng B, 2013. Design-based spatial sampling: theory and implementation. *Environ Model Softw* 40: 280-288. <https://doi.org/10.1016/j.envsoft.2012.09.015>
- Warrick AW, Nielsen DR, 1980. Spatial variability of soil physical properties in the field. In: *Application of soil physics*; Hillel D (ed.). pp: 319-324. Academic Press, NY. <https://doi.org/10.1016/B978-0-12-348580-9.50018-3>