



Harmonization of design-based mapping for spatial populations

A. Marcelli^{1,2} · L. Fattorini³ · S. Franceschi³

Accepted: 21 January 2022
© The Author(s) 2022

Abstract

The mapping of a survey variable throughout a continuum or for finite populations of units is usually performed from a model-dependent perspective. Nevertheless, when a sample of locations/units is selected by a probabilistic sampling scheme, the complex task of modelling can be avoided by using the inverse distance weighting interpolator and deriving the properties of maps in a design-based perspective. Conditions ensuring consistency of maps can be derived mainly based on some obvious assumptions about the pattern of the survey variable throughout the study region as well from the feature of the sampling scheme adopted to select locations/units. Nevertheless, in a design-based setting the totals of the survey variable for a set of domains partitioning the study region are commonly estimated by traditional estimators such as the Horvitz–Thompson estimator in the case of finite populations or the Monte-Carlo estimator in the case of continuous populations or by related estimators exploiting the information of auxiliary variables. That necessarily gives rise to different total estimates with respect to those achieved from the resulting maps as the sum of the interpolated values within domains. To obtain non-discrepant results, a harmonization of maps is here suggested, in such a way that the resulting totals arising from maps coincide with those achieved by traditional estimation. The capacity of the harmonization procedure to maintain consistency is argued theoretically and checked by a simulation study performed on some real populations.

Keywords Inverse distance weighting interpolation · Horvitz–Thompson estimator · Monte Carlo estimator · Simulation study

1 Introduction

Spatial phenomena frequently need detailed information especially in natural resources management. Therefore, mapping becomes crucial for the understanding of spatial patterns of an interest attribute. The population to be surveyed can be a continuum, a finite collection of areas portioning a study region, or a finite collection of points spread throughout a study region. In most cases, the available resources render impractical the complete census

of these populations. Therefore, the survey variable is recorded only for a subset of locations/areas/points and an estimation criterion is adopted to estimate the values of the survey variable within non recorded locations/areas/points and obtain wall-to-wall maps of the interest variable throughout the whole population.

Usually, mapping is approached in a model-based framework: locations/area/points where the variable is recorded are considered as fixed, while values are assumed to be outcomes of a superpopulation probability model (e.g., Cressie 1993). As an alternative approach, Fattorini et al. (2018a, b, 2019) propose mapping continuous populations, finite populations of areas and finite populations of points in a design-based framework. In this scenario, values are viewed as fixed constants and the probability distribution of any sample statistic is determined from the uncertainty entailed by the probabilistic sampling scheme adopted for selecting locations/areas/points. In their works, the authors highlighted as any design-based mapping is challenging. Indeed, when estimating the values of the survey variable for a single location/area/point,

✉ S. Franceschi
sara.franceschi@unisi.it

¹ Department for Innovation in Biological, Agro-Food and Forest Systems, University of Tuscia, Viterbo, Italy

² Department of Sustainable Agro-Ecosystems and Bioresources, Research and Innovation Centre, Fondazione E. Mach, San Michele All'Adige, Italy

³ Department of Economics and Statistics, University of Siena, Siena, Italy

either it has been sampled and therefore there is no need for estimation, or it is unsampled and therefore we do not have any information for performing estimation. Consequently, even in a design-based framework, the use of an assisting model, based on some auxiliary information, seems to be the only way to fill the lack of sample information. As a solution, Fattorini et al. (2018a, b, 2019) considered as assisting model the well-known Tobler's first law of geography, i.e., locations/areas/points close in space tend to have more similar values than those that are far apart (Tobler, 1970). Accordingly, the authors proposed the adoption of the so-called inverse distance weighting interpolator (IDW) (e.g., Henley 1981) and determined the conditions ensuring IDW design-based unbiasedness and consistency.

Nevertheless, it should be acknowledged that frequently, as it occurs for example in forest inventories, estimation of totals and averages of interest attributes for a set of domains partitioning the study region is traditionally performed from a design-based perspective adopting the Horvitz–Thompson (HT) estimator in the case of finite populations or the Monte-Carlo estimator in the case of continuous populations or some modifications able to exploit auxiliary information (e.g., Särndal et al. 1992, chapter 6).

However, total estimates for these domains can be achieved also as totals or integrals of the IDW interpolations within. Obviously, the two resulting sets of estimates will differ. Therefore, the aim of this work is to harmonize IDW maps by rescaling the interpolated values in such a way that totals or integrals of single estimates within domains will match the traditional design-based estimates, avoiding unsuitable discrepancies in the final results, that may be perplexing especially in a report phase. It is worth noting that the harmonization procedure is not introduced to improve the design-based performance of the IDW interpolation, but just to avoid discrepancies of the total estimates achieved from maps with those achieved from traditional methods. Therefore, the capacity of the harmonization procedure to maintain the design-based consistency ensured by the original IDW interpolation has been the main concern. Accordingly, the main target of this paper is to theoretically prove the harmonization consistency and to check it empirically by simulations.

Because averages are achieved as totals divided by the size of the study region in continuous populations or by totals divided by the population size in finite populations of areas or points, and because these quantities must necessarily be known to perform mapping, harmonization with respect to totals or averages are equivalent from a statistical point of view. Moreover, regarding the types of populations considered for harmonization, it should be noticed that for point populations over large study regions

(e.g., plants, shrubs, or trees), the list and locations of population units are not available. As consequence, maps of these populations are precluded. The sole cases in which mapping is possible occur for forest stands located on surfaces of limited size (few hectares) in which 3P sampling can be performed. Indeed, in these cases it becomes possible to visit (and then to locate and list) all the trees by a team of experts and to give predictions of the survey variable for each tree, that is a necessary step to perform 3P sampling (Gregoire and Valentine 2008). However, in most of these cases, the principal aim is the estimation of the total timber volume, while mapping is less urgent: this is due to the small size of the stand and also because the predictions from a crew of experts are likely to well depict the distribution of the timber volumes throughout the stand. For these reasons, harmonization in finite populations of points is not considered in the paper.

The paper is organized as it follows: in Sect. 2 some preliminary results on design-based mapping of continuous populations and finite populations of areas are given. Then, the procedures for harmonizing IDW maps are described in Sect. 3 and 4 where the design-based consistency of the harmonized maps is theoretically proven, and a pseudo-population bootstrap estimator of map precision is proposed. In Sect. 5 a simulation study is performed on a set of real populations to empirically check the capacity of the harmonization procedure in maintaining the consistency of the resulting maps. Concluding remarks are reported in Sect. 6.

2 Preliminary results on design-based IDW interpolation

Consider a study region A , that is supposed to be a connected and compact set of \mathbb{R}^2 , and let f be a bounded measurable function related to the values of a survey variable Y and defined on a subset $B \subset A$. Moreover, let $\|\cdot\|$ be a norm in \mathbb{R}^2 and $\phi: [0, \infty) \rightarrow \mathbb{R}^+$ be a nonincreasing continuous distance function on $(0, \infty)$, with $\phi(0) = 0$ and

$$\lim_{d \rightarrow 0^+} \phi(d) = \infty. \quad (1)$$

A widely applied class of distance functions satisfying (1) is the class of negative powers of order α , given by

$$\phi(d) = d^{-\alpha}, \alpha > 0 \quad (2)$$

(e.g., Gong et al. 2014; Noori et al. 2014; Bărbulescu et al. 2021). The choice of (2) is particularly appealing owing to its simplicity and because of the straightforward interpretation of the α parameter. Indeed, as showed in the next subsections, the IDW interpolator is a convex

combination of the sample values, with weights proportional to the values of the distance function ϕ . Therefore, negative powers of distances obviously give less weight to sample values further away from the location where interpolation is performed with α that plays the role of the smoothing parameter, i.e., for smaller values of α , the interpolated map becomes smoother, till, in the limiting case of $\alpha = 0$, the map is constant as the interpolated values are all equal to the sample mean. On the other hand, when α becomes larger and larger, the estimated map becomes rougher and rougher till, for α approaching infinity, the IDW interpolator reduces to the well-known nearest neighbour (NN) interpolator (Fattorini et al. 2021).

Depending on the features of the spatial populations, the IDW mapping of the interest attribute is performed in the following two settings.

2.1 Continuous populations

B coincides with A and $f(\mathbf{p})$ is the value or the density of the survey variable Y at $\mathbf{p} \in B$. Therefore, mapping necessitates the knowledge of $f(\mathbf{p})$ for each $\mathbf{p} \in B$. To this purpose, let $\mathbf{P}_1, \dots, \mathbf{P}_n$ be n random variables with values in B that represent the n locations selected from B by means of a probabilistic sampling scheme. Then, in accordance with Fattorini et al. (2018a), if the existence of the continuous probability density function of $(\mathbf{P}_1, \dots, \mathbf{P}_n)$ is assumed, for a fixed α , the IDW interpolator of $f(\mathbf{p})$ is almost certainly equal to

$$\hat{f}_\alpha(\mathbf{p}) = \sum_{i=1}^n w_{i,\alpha}(\mathbf{p})f(\mathbf{P}_i), \quad \mathbf{p} \in B \tag{3}$$

with weights that under the distance function (2) turn out to be

$$w_{i,\alpha}(\mathbf{p}) = \frac{\|\mathbf{P}_i - \mathbf{p}\|^{-\alpha}}{\sum_{h=1}^n \|\mathbf{P}_h - \mathbf{p}\|^{-\alpha}}, \quad i = 1, \dots, n.$$

For $\alpha \rightarrow \infty$ the interpolator (3) reduces to the NN interpolator, that is almost certainly equal to

$$\hat{f}_\infty(\mathbf{p}) = f(\mathbf{P}_{NN(p)}) \tag{4}$$

where $\mathbf{P}_{NN(p)} = \operatorname{argmin}_{i=1, \dots, n} \|\mathbf{P}_i - \mathbf{p}\|$. In this case the resulting map is a surface piecewise constant on the Voronoi cells around the sampled locations (Fattorini et al. 2021).

2.2 Finite populations of areas

A is partitioned into a finite population U of N spatial units a_1, \dots, a_N of extents $\lambda_1, \dots, \lambda_N$. In this case B is the set of centroids $\mathbf{b}_1, \dots, \mathbf{b}_N$ of the areas and y_j is the amount of the survey variable Y within a_j . Therefore, mapping

necessitates the knowledge of y_j for each $j \in U$. However, since the area extents are usually known, knowledge of the y_j s is equivalent to the knowledge of densities $f_j = f(\mathbf{b}_j) = y_j/\lambda_j$ for each $j \in U$. Accordingly, if S denotes the set of labels identifying the selected areas, then in accordance with Fattorini et al. (2018b), for a fixed α , the IDW interpolator of f_j is given by

$$\hat{f}_{j,\alpha} = Z_j f_j + (1 - Z_j) \sum_{i \in S} w_{ij,\alpha} f_i, \quad j \in U \tag{5}$$

where Z_j is the random variable equal to 1 if $j \in S$ and equal to 0 otherwise, with weights that under the distance function (2) turn out to be

$$w_{ij,\alpha} = \frac{\|\mathbf{b}_i - \mathbf{b}_j\|^{-\alpha}}{\sum_{h \in S} \|\mathbf{b}_h - \mathbf{b}_j\|^{-\alpha}}, \quad i \in S.$$

For $\alpha \rightarrow \infty$ the interpolator (5) reduces to the NN interpolator

$$\hat{f}_{j,\infty} = Z_j f_j + \frac{(1 - Z_j)}{\operatorname{Card}(H_j)} \sum_{i \in H_j} f_i, \quad j \in U \tag{6}$$

where $H_j = \{i : \|\mathbf{b}_i - \mathbf{b}_j\| = \min_{i \in S} \|\mathbf{b}_h - \mathbf{b}_j\|\}$ is the set of centroids in the sample that are nearest to \mathbf{b}_j . Indeed, contrary to the case of continuous populations, in the case of finite populations of areas the nearest neighbours of an area may be more than 1 as, for example, in the case of populations of regular polygons such as pixels (Fattorini et al. 2021).

Once the $\hat{f}_{j,\alpha}$ s are achieved, the interpolation of the y_j s is simply given by $\hat{y}_{j,\alpha} = \lambda_j \hat{f}_{j,\alpha}$ for $j \in U$. However, the interpolation of densities is more suitable for working in the asymptotic scenario considered by Fattorini et al. (2018b, 2021) in which the λ_j s decrease and then the y_j s approach zero.

2.3 Consistency conditions

The asymptotic properties of the IDW interpolators (3) and (5) are derived in Fattorini et al. (2018a, b), respectively, while those of their NN counterparts (4) and (6), achieved for $\alpha \rightarrow \infty$, are derived in Fattorini et al. (2021). Without going into the technical details provided in those papers, for both the population settings, the design-based unbiasedness and consistency of IDW and NN interpolators concern: i) a sort of smoothness of the function f onto the study region with jumps that occur in sub-sets of measure 0; ii) when dealing with finite populations of areas, the regularities in the area shapes; iii) the capacity of the sampling design to asymptotically achieve a spatial balance of the selected locations/areas; iv) the use of negative power distance functions with $\alpha > 2$.

It is worth noting that the four conditions seem to match most real situations encountered in environmental surveys. Indeed, the smoothness assumption i) is very common and it is at the basis of most interpolation techniques (e.g., Cressie 1993, Sect. 3.1). Moreover, assumption i) is also reasonably valid in many natural scenarios where the density of an attribute changes smoothly throughout space (continuity) and when it changes abruptly, that usually occurs along borders delineating variations in the characteristics of the study region (e.g., forest-meadows). Therefore, borders may be realistically approximated by curves well approaching the theoretical condition of discontinuity over a region of zero measure. Regarding assumption ii) concerning the regularity of the shape of areas, it is ensured by the fact that in most cases, especially in forest and vegetation surveys, areas are regular polygons. Finally, iii) and iv) actually do not constitute assumptions because both sampling schemes and distance functions are chosen by the user. In particular, as emphasized in Sect. 3, the asymptotic spatial balance required by iii) is ensured under the schemes usually applied in environmental surveys.

2.4 Data-driven choice of α

Because any value of $\alpha > 2$, including $\alpha = \infty$, ensures the design-based consistency of the IDW interpolation, Fattorini et al. (submitted) propose to choose α , that in this framework plays the important role of a smoothing parameter, by means of a data driven procedure. In particular, the authors propose the use of the leave-one-out cross validation (LOOCV) that constitutes an intuitive and widely applied technique in spatial interpolation (e.g., Giraldo et al. 2011; Ignaccolo et al. 2014; Montanari and Cicchitelli 2014).

LOOCV consists in removing one location/area at a time from the sampled ones, interpolating the value or density of the survey variable at the removed location/area using all other locations/areas in the sample and then repeating this process for each location/area in the sample. The interpolated values at each sample location/area are then compared with the actual values minimizing the sum of squared differences.

Accordingly, in the case of continuous populations, α is selected to minimize

$$SSD(\alpha) = \sum_{i=1}^n \left[\hat{f}_{-i,\alpha}(\mathbf{P}_i) - f(\mathbf{P}_i) \right]^2 \tag{7}$$

where $\hat{f}_{-i,\alpha}(\mathbf{P}_i)$ is the IDW interpolator of $f(\mathbf{P}_i)$ achieved by means of (3) or (4) from the sample of $n - 1$ locations obtained deleting the sample location \mathbf{P}_i .

Similarly, in the case of finite populations of areas, α is selected to minimize

$$SSD(\alpha) = \sum_{i \in S} (\hat{f}_{-i,\alpha} - f_i)^2 \tag{8}$$

where $\hat{f}_{-i,\alpha}$ is the IDW interpolator of f_i achieved by means of (5) or (6) from the sample of $n - 1$ areas obtained deleting the area $i \in S$. In the following, the IDW interpolators with α chosen by means of LOOCV are denoted by $\hat{f}_{\hat{\alpha}}(\mathbf{p})$ for each $\mathbf{p} \in B$ in the case of continuous populations or by $\hat{f}_{j,\hat{\alpha}}$ for each $j \in U$ in the case of finite populations of areas and are referred to as data-driven IDW (DD-IDW) interpolators.

Because $\hat{\alpha}$ is chosen from sample data instead of being fixed in advance, it is a random variable whose presence may, in principle, increase the variability of the DD-IDW interpolators and may preclude their consistency. Fattorini et al. (submitted) broad the consistency results achieved for IDW interpolators proving the consistency of the DD-IDW under the same asymptotic scenarios.

2.5 Bootstrap estimation of precision

Regarding the estimation of the precision of the DD-IDW interpolators, Fattorini et al. (submitted) propose the use of a pseudo-population bootstrap (see e.g., Quatemberg 2015) based on adopting the DD-IDW map achieved from the sample as the pseudo-population from which M bootstrap samples are re-sampled using the same scheme adopted to select the original sample and then achieving as many DD-IDW bootstrap maps from these samples.

Accordingly, in the case of continuous populations, let $\hat{f}_{\hat{\alpha}}(B) = \{ \hat{f}_{\hat{\alpha}}(\mathbf{p}), \mathbf{p} \in B \}$ be the DD-IDW map based on the sample values $f(\mathbf{P}_1), \dots, f(\mathbf{P}_n)$. Then, for each $\mathbf{p} \in B$, the pseudo-population bootstrap estimator of the root mean squared error (RMSE) of $\hat{f}_{\hat{\alpha}}(\mathbf{p})$ is given by

$$rmse_M^*(\mathbf{p}) = \left\{ \frac{1}{M} \sum_{m=1}^M \left[\hat{f}_{\hat{\alpha}_m^*}(\mathbf{p}) - \hat{f}_{\hat{\alpha}}(\mathbf{p}) \right]^2 \right\}^{\frac{1}{2}}, \mathbf{p} \in B \tag{9}$$

where $\mathbf{P}_{1,m}^*, \dots, \mathbf{P}_{n,m}^*$ are the locations selected in the m -th bootstrap resampling by means of the scheme adopted to select the original sample, $\hat{f}_{\hat{\alpha}}(\mathbf{P}_{1,m}^*), \dots, \hat{f}_{\hat{\alpha}}(\mathbf{P}_{n,m}^*)$ are the values at these locations derived from the estimated map $\hat{f}_{\hat{\alpha}}(B)$, $\hat{\alpha}_m^*$ is the LOOCV choice of α performed on the bootstrap sample and $\hat{f}_{\hat{\alpha}_m^*}(\mathbf{p})$ is the DD-IDW interpolation at $\mathbf{p} \in B$ based on $\hat{\alpha}_m^*$.

Similarly, in the case of finite populations of areas, let $\hat{f}_{\hat{\alpha}}(U) = \{ \hat{f}_{j,\hat{\alpha}}, j \in U \}$ be the DD-IDW map based on the sample densities $f_j, j \in S$. Then, for each $j \in U$, the pseudo-

population bootstrap estimator of the RMSE of $\hat{f}_{j,\hat{\alpha}}$ is given by

$$\widehat{rmse}_{j,M}^* = \left\{ \frac{1}{M} \sum_{m=1}^M \left(\hat{f}_{j,\hat{\alpha}_m}^* - \hat{f}_{j,\hat{\alpha}} \right)^2 \right\}^{\frac{1}{2}}, j \in U \tag{10}$$

where S_m^* is the sample of areas selected in the m -th bootstrap resampling by means of the scheme adopted to select the original sample, $\hat{f}_{j,\hat{\alpha}_m}^*, j \in S_m^*$ are the densities in these areas derived from the estimated map $\hat{f}_{\hat{\alpha}}(U)$, $\hat{\alpha}_m^*$ is the LOOCV choice of α performed on the bootstrap sample and $\hat{f}_{j,\hat{\alpha}_m}^*$ is the DD-IDW interpolation of $f_j, j \in U$, based on $\hat{\alpha}_m^*$.

Fattorini et al. (submitted) prove that for large sample sizes and for a sufficiently large M , the bootstrap estimators (9) and (10) tend to be conservative with the ratio of the expectation of the bootstrap RMSE to the true value that is bounded by $\sqrt{10}$. Even if this result may induce to suspect a large overestimation of the RMSEs, such bound should be viewed as a threshold limiting possible overestimation.

3 Harmonization with the overall population total

In environmental and ecological studies, the total of a survey variable for the whole study region frequently covers a role of great interest. As typical examples, estimation of the total amount of a pollutant in a lake is crucial for achieving information on environmental changes in the surrounding zones (e.g., Greaver et al. 2016), estimation of the total erosion extent in a region is fundamental for the management of cultivations (e.g., Kelley, 1990), estimation of the total carbon storage in a forest is essential for determining sequestration capacity (e.g., FAO 2010, Chapter 2).

Once locations or areas are selected by means of a probabilistic sampling scheme, the population total is commonly estimated by means of the HT criterion or by related criteria able to exploit the presence of auxiliary information. These criteria constitute consolidated, widely applied strategies as they were unbiased or approximately unbiased with design-based variances with known analytic expressions or approximations. Moreover, variances can be estimated on the basis of those expressions/approximations (Gregoire and Valentine, 2008).

However, total estimates naturally arise also from the resulting maps as the integral, in the continuous case, or the sum, in the discrete cases, of the interpolated values, thus achieving total estimates that invariably differ from those achieved by traditional, HT-based techniques. To eliminate

these unsuitable discrepancies, the matching of the two total estimates can be performed by rescaling the DD-IDW interpolations of the interest attribute for each location/area. Accordingly, for continuous populations, the rescaled DD-IDW map is given by

$$\tilde{f}_{\hat{\alpha}}(\mathbf{p}) = \frac{\hat{T}}{\hat{T}_{\hat{\alpha}}} \hat{f}_{\hat{\alpha}}(\mathbf{p}), \mathbf{p} \in B \tag{11}$$

while, for finite populations of areas, the rescaled DD-IDW map is given by

$$\tilde{f}_{j,\hat{\alpha}}^* = \frac{\hat{T}}{\hat{T}_{\hat{\alpha}}} \hat{f}_{j,\hat{\alpha}}^*, j \in U \tag{12}$$

respectively, where, in both cases, \hat{T} denotes the total estimate obtained by means of a commonly adopted HT-based technique and $\hat{T}_{\hat{\alpha}}$ denotes the total estimate obtained from the resulting DD-IDW map.

3.1 Harmonization for continuous populations

In the case of continuous populations, the spatial total can be expressed as

$$T = \int_B f(\mathbf{p}) d\mathbf{p} \tag{13}$$

(see e.g., Stevens 1997; Gregoire and Valentine 2008, Chapter 10), while its design-based estimation \hat{T} can be obtained extending the HT estimator to the continuous case

$$\hat{T} = \sum_{i=1}^n \frac{f(\mathbf{P}_i)}{\pi(\mathbf{P}_i)} \tag{14}$$

where $\pi(\mathbf{P}_i)$ denotes the strictly positive inclusion density function on B at the sample locations $\mathbf{P}_i (i = 1, \dots, n)$ (see e.g., Cordy, 1993).

Because of the integral representation (13), the estimation of T may be approached as a Monte Carlo integration. Interestingly, the most common Monte Carlo integration methods such as crude Monte Carlo integration, modified Monte Carlo integration and random grid Monte Carlo integration are equivalent to Uniform Random Sampling (URS), Tessellation Stratified Sampling (TSS) and Systematic Grid Sampling (SGS), respectively, that constitute the most common sampling schemes adopted in environmental surveys for continuous populations (e.g., Barabesi 2003). In these cases, the estimator (14) reduces to

$$\hat{T} = \frac{\lambda(B)}{n} \sum_{i=1}^n f(\mathbf{P}_i). \tag{15}$$

Consistency of (15) under URS, TSS and SGS has been proven by supposing a sequence of designs, each of them

selecting an increasing number of locations. In particular, the relative efficiency of TSS with respect to URS has been proven for finite samples, also proving that efficiency approaches infinity as the sample size increases because TSS variance goes to 0 more quickly than under URS (Barabesi and Franceschi 2011; Barabesi et al. 2012, 2015). Moreover, consistency of (15) under SGS has been proven by Fattorini et al. (2020).

On the other hand, from the resulting DD-IDW map, the total estimate turns out to be

$$\widehat{T}_\alpha = \int_B \widehat{f}_\alpha(\mathbf{p}) d\mathbf{p} = \sum_{i=1}^n f(\mathbf{P}_i) \int_B w_{i,\alpha}(\mathbf{p}) d\mathbf{p}. \tag{16}$$

For α fixed and under the same asymptotic scenario and the same schemes (URS, TSS and SGS) that ensure the design-based consistency of (15), Fattorini et al. (2018a) prove the design-based consistency of the IDW estimator $\widehat{f}_\alpha(\mathbf{p})$ that, in turn,—owing to the Lebesgue dominated convergence Theorem — entails the consistency of the integral $\widehat{T}_\alpha = \int_B \widehat{f}_\alpha(\mathbf{p}) d\mathbf{p}$ to the true total T . Therefore, stated the consistency of the DD-IDW interpolator $\widehat{f}_\alpha(\mathbf{p})$ under the same asymptotic scenario, that also entails the consistency of (16).

Joining the two consistency results for the estimators (15) and (16), the rescaling constant $\widehat{T}/\widehat{T}_\alpha$ in Eq. (11) obviously converges to 1, in such a way that the harmonized interpolator $\widehat{f}_\alpha(\mathbf{p})$ converges to the DD-IDW interpolator $\widehat{f}_\alpha(\mathbf{p})$. That proves consistency of the harmonized interpolator (11) for continuous populations under URS, TSS and SGS.

3.2 Harmonization for finite populations of areas

In the case of finite populations of areas, the population total can be expressed as

$$T = \sum_{j \in U} y_j \tag{17}$$

while its design-based estimation can be obtained by means of the HT estimator

$$\widehat{T} = \sum_{j \in S} \frac{y_j}{\pi_j} \tag{18}$$

where π_j denotes the first-order inclusion probability of the area j induced by the sampling scheme adopted to select areas. As the extents of areas partitioning the study region decrease in such a way that their number and sample size increase, Fattorini et al. (2020) derived conditions ensuring design-based consistency of (18). In particular,

they proved that consistency conditions are satisfied when Simple Random Sampling Without Replacement (SRSWOR) and One Per Stratum Sampling (OPSS), are adopted. Moreover, they proved consistency also under Systematic Sampling (SYS), that however necessitates further assumptions.

On the other hand, from the resulting DD-IDW map, the total estimate turns out to be

$$\widehat{T}_\alpha = \sum_{j \in U} \lambda_j \widehat{f}_{j,\alpha}. \tag{19}$$

For α fixed and the same asymptotic scenario and the same schemes (SRSWOR, OPSS and SYS) that ensure consistency of (18), Fattorini et al. (2018b) prove the design-based consistency of the IDW estimator $\widehat{f}_{j,\alpha}$ that, in turn—owing to the Result B.1, Appendix B of that paper — entails the consistency of $\widehat{T}_\alpha = \sum_{j \in U} \lambda_j \widehat{f}_{j,\alpha}$ to the true total T . Therefore, stated the consistency of the DD-IDW interpolator $\widehat{f}_{j,\alpha}$ under the same asymptotic scenario, that also entails the consistency of (19).

Joining the two consistency results for the estimators (18) and (19), the rescaling constant $\widehat{T}/\widehat{T}_\alpha$ in Eq. (12) obviously converges to 1, in such a way that the harmonized interpolator $\widehat{f}_{j,\alpha}$ converges to the DD-IDW interpolator $\widehat{f}_{j,\alpha}$. That proves consistency of the harmonized interpolator (12) for finite populations of areas under SRSWOR, OPSS and SYS.

3.3 Bootstrap estimation of precision

Regarding the estimation of the precision of the harmonized interpolators, because harmonization is performed from sample data, rescaling the DD-IDW interpolations by means of the ratio $\widehat{T}/\widehat{T}_\alpha$ involves uncertainty that must be accounted in the bootstrap procedure. Therefore, in the case of continuous population, the bootstrap RMSE estimator (9) is changed into

$$rmse_M^*(\mathbf{p}) = \left\{ \frac{1}{M} \sum_{m=1}^M \left[\widehat{f}_{\alpha_m}^*(\mathbf{p}) - \widehat{f}_\alpha(\mathbf{p}) \right]^2 \right\}^{\frac{1}{2}}, \mathbf{p} \in B \tag{20}$$

where $\widehat{f}_{\alpha_m}^*(\mathbf{p})$ is the harmonized counterparts of $\widehat{f}_{\alpha_m}^*(\mathbf{p})$ achieved rescaling $\widehat{f}_{\alpha_m}^*(\mathbf{p})$ by the ratio $\widehat{T}^*/\widehat{T}_{\alpha_m}^*$ where \widehat{T}^* and $\widehat{T}_{\alpha_m}^*$ are the total estimates (15) and (16) achieved from the bootstrap sample values $\widehat{f}_\alpha^*(\mathbf{P}_{1,m}^*), \dots, \widehat{f}_\alpha^*(\mathbf{P}_{n,m}^*)$.

In the case of finite populations of areas, the bootstrap RMSE estimator (10) is changed into

$$rm\hat{se}_{j,M}^* = \left\{ \frac{1}{M} \sum_{m=1}^M \left(f_{j,\hat{\alpha}_m}^* - \hat{f}_{j,\hat{\alpha}} \right)^2 \right\}^{\frac{1}{2}}, j \in U \tag{21}$$

where $f_{j,\hat{\alpha}_m}^*$ is the harmonized counterparts of $\hat{f}_{j,\hat{\alpha}_m}^*$ achieved rescaling $\hat{f}_{j,\hat{\alpha}_m}^*$ by the ratio $\hat{T}^* / \hat{T}_{\hat{\alpha}_m}^*$ where \hat{T}^* and $\hat{T}_{\hat{\alpha}_m}^*$ are the total estimates (18) and (19) achieved from the bootstrap sample values $\hat{f}_{j,\hat{\alpha}}, j \in S_m^*$.

4 Harmonization by domains

It frequently occurs that total estimates of an interest attribute are also required for D subpopulations partitioning the entire population, usually referred to as domains (e.g., Särndal et al. 1992 Sect. 10.3). For example, a survey region may be divided into D domains by administrative bounds (e.g., regions, counties, municipalities) and we could be interested, besides to the mapping and to the overall total of an attribute, to its sub-totals within the domains.

In the case of continuous populations, denote by B_1, \dots, B_D the D subsets partitioning B , whose totals are of interest. In this case, estimation of totals within domains can be performed simply defining, for each domain $d = 1, \dots, D$, the function.

$$f_{(d)}(\mathbf{p}) = \begin{cases} f(\mathbf{p}) & \text{if } p \in B_d \\ 0 & \text{otherwise} \end{cases}$$

in such a way that the total for the domain d , say $T_{(d)}$, is obtained as in (13) by substituting $f(\mathbf{p})$ with $f_{(d)}(\mathbf{p})$. Similarly, the Monte Carlo estimator for the d -th domain, say $\hat{T}_{(d)}$, can be performed as in (15), once again substituting $f(\mathbf{p})$ with $f_{(d)}(\mathbf{p})$. On the other hand, from the resulting DD-IDW map, the total estimate for each domain $d = 1, \dots, D$, say $\hat{T}_{\hat{\alpha}(d)}$, turns out to be as in (16) with the integral extended to B_d , instead of the whole B . Therefore, for continuous populations, the rescaled map that ensures the matching of total estimates for each domain as well as the matching of the overall total estimates is given by

$$\tilde{f}_{\hat{\alpha}}(\mathbf{p}) = \frac{\hat{T}_{(d)}}{\hat{T}_{\hat{\alpha}(d)}} \hat{f}_{\hat{\alpha}}(\mathbf{p}), \mathbf{p} \in B_d, d = 1, \dots, D. \tag{22}$$

Consistency of (22) arises, *mutatis mutandis*, from the same consideration performed in Sect. 3.1.

In the case of finite populations of areas, denote by U_1, \dots, U_D the D subsets of areas partitioning U , whose totals are of interest. In this case, estimation of totals within domains can be performed simply defining, for each domain $d = 1, \dots, D$, the values.

$$f_{j(d)} = \begin{cases} f_j & \text{if } j \in U_d \\ 0 & \text{otherwise} \end{cases}$$

in such a way that the total for the domain d , say $T_{(d)}$, is obtained as in (17) by substituting f_j with $f_{j(d)}$. Similarly, the HT estimator for the d -th domain, say $\hat{T}_{(d)}$, can be performed as in (18), once again substituting f_j with $f_{j(d)}$. On the other hand, from the resulting DD-IDW map, the total estimate for each domain $d = 1, \dots, D$, say $\hat{T}_{\hat{\alpha}(d)}$, turns out to be as in (19) with the summand extended to U_d , instead of the whole U . Therefore, for finite populations of areas, the rescaled map ensuring the matching of total estimates for each domain as well as the matching of the overall total estimates is given by

$$\tilde{f}_{j,\hat{\alpha}} = \frac{\hat{T}_{(d)}}{\hat{T}_{\hat{\alpha}(d)}} \hat{f}_{j,\hat{\alpha}}, j \in U_d, d = 1, \dots, D. \tag{23}$$

Consistency of (23) arises, *mutatis mutandis*, from the same consideration performed in Sect. 3.2.

Finally, regarding the bootstrap RMSE estimation, in the case of continuous populations it is performed by means of Eq. (20) substituting $f(\mathbf{p})$ with $f_{(d)}(\mathbf{p})$, and in the case of finite populations of areas it is performed by means of Eq. (21) substituting f_j with $f_{j(d)}$.

5 Simulation studies

For each population setting, a simulation study is performed to empirically check if and how much the harmonization procedure deteriorates the performance of DD-IDW maps, as well as to check the rate of convergence of the harmonized maps to the original ones.

5.1 Populations and sampling

As to continuous populations, the survey region considered for the simulation was a quadrat region of 90,000 ha located in North-Western Tuscany (Central Italy). The population values to be mapped were the precipitations (mm) occurred between 3rd January 2021 and 3rd February 2021, that were artificially achieved by means of an ordinary kriging prediction performed from the values recorded on 32 rain gauge stations that were present in the survey region. The population average was $\bar{Y} = 299$ mm for the whole region (see Fig. 1). Moreover, the region was partitioned into $D = 2, 4, 8$ domains of equal sizes, as depicted in Fig. 2. Precipitation averages (in mm) within the two domains were $\bar{Y}_{(1)} = 306$ and $\bar{Y}_{(2)} = 293$, precipitation averages within the four domains were $\bar{Y}_{(1)} = 255$, $\bar{Y}_{(2)} = 361$, $\bar{Y}_{(3)} = 323$, and $\bar{Y}_{(4)} = 257$, precipitation

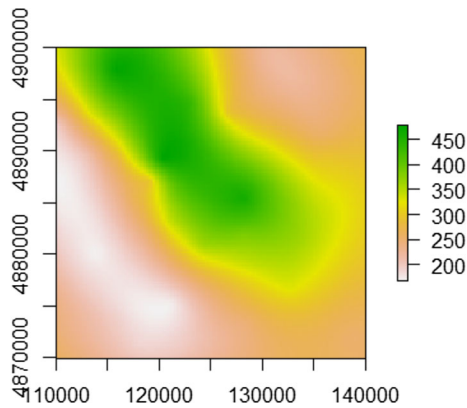


Fig. 1 Precipitations (mm) between 3rd January 2021 and 3rd February 2021 in a quadrat region of North-Western Tuscany (Central Italy)

averages within the eight domains were $\bar{Y}_{(1)} = 219$, $\bar{Y}_{(2)} = 227$, $\bar{Y}_{(3)} = 300$, $\bar{Y}_{(4)} = 423$, $\bar{Y}_{(5)} = 282$, $\bar{Y}_{(6)} = 319$, $\bar{Y}_{(7)} = 341$, and $\bar{Y}_{(8)} = 281$.

Sampling was performed by selecting $n = 16, 36, 64, 100$ locations on the quadrat region by means of URS, TSS and SGS, i.e., selecting n locations completely at random (URS), partitioning the quadrat into n sub-quadrats of equal size and selecting a location at random within each of them (TSS), or selecting one location at random in one of them and then repeating the location in the others (SGS).

As to finite populations of areas, the survey region considered for the simulation was a rectangle of about 212 ha located in Calabria (Southern Italy). Three populations of $N = 250, 1000, 4000$ areas were considered by

Fig. 2 Partition of the region in Fig. 1 into $D = 2, 4, 8$ domains of equal size

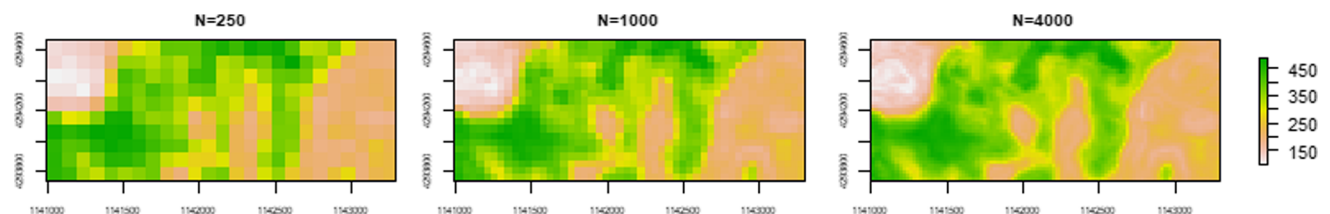
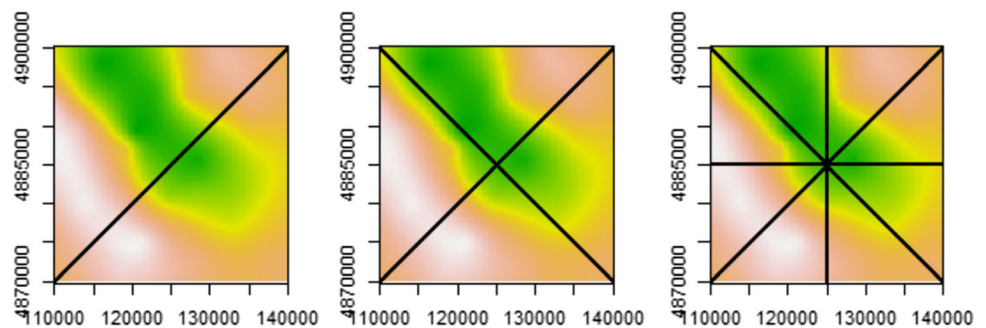


Fig. 3 Maps of the growing stock volumes in a rectangular region located in Calabria (Southern Italy) partitioned into three populations of 250, 1,000 and 4,000 areas

partitioning the region into as many rectangles of size 8464, 2116 and 529 m², respectively. The population values to be mapped were the growing stock volumes (m³/ha) within the areas that were artificially achieved by means of a random forest imputation technique (see Chirici et al. 2020 for details) from the ground data recorded in the 2000–2007 during the Italian National Forest Inventory (Fattorini et al. 2006). The population total resulted $T = 63,524.6$ m³ for the whole region (see Fig. 3). Moreover, the region was partitioned into $D = 2, 4, 8$ domains constituted by an equal number of areas, as depicted in Fig. 4. Totals (in m³) within the two domains were $T_{(1)} = 34,915.5$ and $T_{(2)} = 28,609.1$, totals within the four domains were $T_{(1)} = 16,830.6$, $T_{(2)} = 18,048.9$, $T_{(3)} = 12,634.2$, and $T_{(4)} = 15,974.8$, totals within the eight domains were $T_{(1)} = 8891.6$, $T_{(2)} = 10,553.7$, $T_{(3)} = 6276.9$, $T_{(4)} = 8755.0$, $T_{(5)} = 9329.9$, $T_{(6)} = 6350.5$, $T_{(7)} = 6283.8$, and $T_{(8)} = 7083.3$.

Sampling was performed by selecting samples of $n = N/10$ areas by means of SRSWOR, OPSS and SYS, i.e., selecting n areas at random without replacement (SRSWOR), partitioning the populations into n blocks of 2×5 contiguous areas and selecting an area at random within each of them (OPSS), or selecting one area at random in one block and then repeating it in the others (SYS).

5.2 Simulation

For each combination of population, sampling scheme and sample size, sampling was replicated $R = 10,000$ times. At each simulation run, the DD-IDW map was obtained by

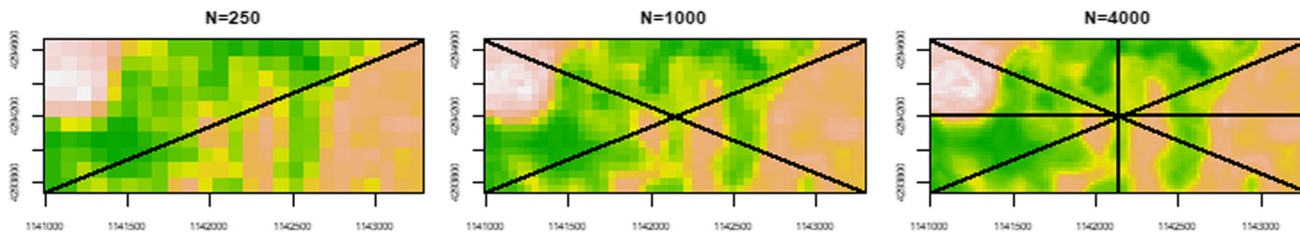


Fig. 4 Partition of the rectangular region in Fig. 3 into $D = 2, 4, 8$ domains constituted by an equal number of areas

means of the LOOCV selection of α , that was performed by minimizing (7) or (8), for the continuous population or finite populations of areas, respectively, starting with $\alpha = 3$ and increasing α by one. Since for quite large values of α the IDW interpolator is practically indistinguishable from the NN interpolator, if the minimum of (7) or (8) was reached for $\alpha = 21$, the NN interpolator was used. Moreover, the bootstrap RMSE estimators (9) or (10) were adopted in the case of the continuous population and finite populations of areas, respectively, computed by means of $M = 1000$ bootstrap samples. Then, at each simulation run, harmonization was performed by estimating the population totals by means of the traditional estimators (15) or (18) in the case of the continuous populations or finite populations of areas, respectively, by estimating the totals from the resulting maps by means of estimators (16) or (19) in accordance with the two types of populations, and then harmonizing the DD-IDW by rescaling them by means of Eq. (11) for the continuous population and by means of Eq. (12) for finite populations of areas. Finally, domain partitions were considered for each population as depicted in Figs. 2 and 4 and harmonization was performed for each of the $D = 2, 4, 8$ domains as described in Sect. 4. Regarding the estimation of precision for the harmonized maps, the bootstrap RMSE estimators (20) or (21) were adopted for the continuous population or the finite populations of areas, respectively, by means of $M = 1000$ bootstrap samples.

5.3 Performance indicators

In the case of the continuous population, interpolation was performed on a regular grid of 100×100 within the quadrat region of Fig. 1. Let $\hat{f}_r(\mathbf{p}_{k,l})$ and $\tilde{f}_r(\mathbf{p}_{k,l})$ be the DD-IDW and the harmonized interpolations of the continuous population computed at the node $\mathbf{p}_{k,l}$ of the grid ($k, l = 1, \dots, 100$) at the r simulation run, and let $\widehat{rmse}_r^*(\mathbf{p}_{k,l})$ and $\widetilde{rmse}_r^*(\mathbf{p}_{k,l})$ be the bootstrap RMSE estimates from Eqs. (9) and (20), respectively. Based on the $R = 10,000$ runs, the absolute bias (AB)

$$AB_{k,l} = \left| \frac{1}{R} \sum_{r=1}^R \hat{f}_r(\mathbf{p}_{k,l}) - f(\mathbf{p}_{k,l}) \right| \tag{24}$$

the RMSE

$$RMSE_{k,l} = \left\{ \frac{1}{R} \sum_{r=1}^R \left[\hat{f}_r(\mathbf{p}_{k,l}) - f(\mathbf{p}_{k,l}) \right]^2 \right\}^{1/2} \tag{25}$$

and the bootstrap ratio (BORAT)

$$BORAT_{k,l} = \frac{\frac{1}{R} \sum_{r=1}^R \widehat{rmse}_r^*(\mathbf{p}_{k,l})}{RMSE_{k,l}} \tag{26}$$

were computed for the DD-IDW interpolation at each node $\mathbf{p}_{k,l}$ for $k, l = 1, \dots, 100$. The same indicators were computed for the harmonized interpolation replacing $\hat{f}_r(\mathbf{p}_{k,l})$ with $\tilde{f}_r(\mathbf{p}_{k,l})$ in Eqs. (24) and (25) and replacing $\widehat{rmse}_r^*(\mathbf{p}_{k,l})$ with $\widetilde{rmse}_r^*(\mathbf{p}_{k,l})$ in Eq. (26).

Similarly, in the case of finite populations of areas, let $\hat{f}_{j,r}$ and $\tilde{f}_{j,r}$ be the DD-IDW and the harmonized interpolations for the area j in the finite populations of areas at the r simulation run, and let $\widehat{rmse}_{j,r}^*$ and $\widetilde{rmse}_{j,r}^*$ be the bootstrap RMSE estimates from Eqs. (10) and (21), respectively. Based on the $R = 10,000$ runs, the AB

$$AB_j = \left| \frac{1}{R} \sum_{r=1}^R \hat{y}_{j,r} - y_j \right| \tag{27}$$

the RMSE

$$RMSE_j = \left\{ \frac{1}{R} \sum_{r=1}^R (\hat{y}_{j,r} - y_j)^2 \right\}^{1/2} \tag{28}$$

and the BORAT

$$BORAT_j = \frac{\frac{1}{R} \sum_{r=1}^R \widehat{rmse}_{j,r}^*}{RMSE_j} \tag{29}$$

were computed for each area j . The same indicators were computed for the harmonized interpolation replacing $\hat{f}_{j,r}$ with $\tilde{f}_{j,r}$ in Eqs. (27) and (28) and replacing $\widehat{rmse}_{j,r}^*$ with $\widetilde{rmse}_{j,r}^*$ in Eq. (29).

Finally, let \widehat{T}_r and $\widehat{T}_{\alpha,r}$ be the total estimates (overall or within a domain) achieved at the r simulation run by means of the HT or Monte Carlo estimators and from the DD-IDW maps, respectively. The Monte Carlo distribution of the correction factors

$$RATIO_r = \frac{\widehat{T}_r}{\widehat{T}_{\alpha,r}} \quad r = 1, \dots, R \quad (30)$$

was considered to evidence the level of matching between the DD-IDW and the harmonized maps.

5.4 Results

Tables SI1 and SI6 in the Online Resource file contain the minima, the averages, and the maxima of ABs, RMSEs, and BORATs arising from the DD-IDW interpolation for each population, sampling scheme and sample size. Figures SI1-SI6 in the Online Resource file show the spatial patterns of these performance indicators. The same indicators are reported in Tables SI2 and SI7 for the harmonized interpolation with respect to the overall total. Figures showing the spatial pattern of these indicators are not reported because they resulted quite similar to those achieved using the DD-IDW interpolation. Moreover, for each population, sampling scheme and sample size, Tables SI3 and SI8 report minima, averages, and maxima of the Monte Carlo distributions of the correction factor (30). Finally, Tables SI4 and SI9 contain the minima, the averages, and the maxima of ABs, RMSEs, and BORATs arising when harmonization is performed with respect to the total estimates within $D = 2, 4, 8$ domains, while Tables SI5 and SI10 report minima, averages, and maxima of the pooled Monte Carlo distributions of the $R \times D$ correction factors (30) in presence of $D = 2, 4, 8$ domains.

Simulation results show that harmonized and DD-IDW maps are comparable in terms of ABs, RMSEs and BORATs for all the populations and sampling schemes, suggesting that harmonization can be performed without relevant loss in accuracy, precision, and bootstrap performance. Furthermore, in all the situations, the correction factor (30) quickly approaches to 1 as the sample sizes increase. However, the RMSEs of the harmonized maps tend to increase as the number of domains increases. This is an expected result due to the fact that, when the number of domains is large and their sizes decrease, the number of selected units within domains becomes smaller and smaller, thus invariably reducing the precision of the total estimates from which harmonization is performed.

6 Final remarks

The recent papers by Fattorini (2018a, 2018b, 2019) propose a novel approach for mapping continuous populations or finite populations of areas and points in a design-based framework, avoiding the complex task of modelling the surfaces or the finite populations to be mapped. However, in the design-based approach, an unresolved problem takes rise. Because totals and averages throughout the whole survey region or within domains are traditionally achieved by the HT or related criteria, these estimates invariably differ from the estimates achieved by the resulting maps. For avoiding these unsuitable discrepancies, that would appear awkward especially in a reporting phase, we here propose to rescale the resulting maps in such a way that the total or average estimates arising from maps will match the traditional estimates. We also prove the asymptotic convergence of the two maps, i.e., the convergence to one of the rescaling factors. That is unambiguously confirmed by the simulation study, as the Monte Carlo distributions of the rescaling factors quickly approach one as sizes increase, even for moderately small domains. Therefore, if we simply look at those results, we may argue that harmonization is a useless procedure because DD-IDW and harmonized maps are very similar. However, if we look at the minima and maxima of the rescaling factor, reported in Online Resource file, it is apparent that in some situations, especially when the estimation of domain totals is involved, the rescaling constant may vary from about 0.1 to 3.4.

On the other hand, the deterioration of precision of harmonized maps when domain sizes decrease ought to warn against an acritical use of harmonization. Obviously, as the number of domains increases and their extents become small, the number of sample units available to perform estimation within domains becomes small and, in some cases, may be 0. This fact will necessarily introduce an increase in the variance in the estimator based on the HT or related criteria. As we cannot trust in these estimates, harmonization with respect to small domains should be avoided. These issues are well known in literature as small area estimation problems, for which an extensive body of knowledge has emerged recently (see e.g., Rao and Molina, 2015 and references therein).

At the end of this paper it should be pointed out that the harmonization problem could be completely bypassed if a model-based mapping was adopted. Indeed, in presence of spatial autocorrelation and second-order stationarity of the spatial process that is supposed to generate the population under study, the kriging methods not only provide the best linear unbiased mapping but directly provides the best linear predictor of totals within subareas with no necessity

of harmonization (e.g. Thompson, 2012, Sect. 20.4). In addition, a long sequence of alternative model-based mapping procedures is available, such as sandwich mapping (Wang et al. 2013) that is suitable in presence of weak spatial autocorrelation and spatial heterogeneity. Indeed, when dealing with some populations, past experience may have established convincingly that certain types of patterns are typical for the survey variable. In these cases, as pointed out by Wang et al. (2012) the patterns can be used to identify a family of stochastic models (superpopulations) that are supposed to generate the population under study for obtaining the most precise possible mapping for a given amount of sampling effort. Obviously, also in design-based approach, in which the minimal sufficient statistics is the unordered set of distinct labelled observations (e.g. Basu 1969), one would like to be able to say what estimator produces the best mapping. However, because in this case the minimal sufficient statistic is not complete (Thompson, 2012, Sect. 9.4), one cannot make statements about one mapping method being best. The lack of completeness of the minimal sufficient statistic in design-based inference has been the main reason for lack of optimality results in this approach. Therefore, one may wonder why to adopt a design-based mapping with the subsequent necessity of performing harmonization if that problem disappears in model-based approaches with also the possibility of achieving the best mapping.

The contraposition between model-based and design-based approaches is a deeply debated issue. Drawbacks and merits of the two approaches are well delineated in both statistical literature (e.g., de Gruijter and ter Braak, 1990; Smith 1994, 2001; Little 2004; Thompson, 2012) as well as in environmental applications (e.g., Schreuder et al. 1993; Gregoire, 1998; Gregoire and Valentine, 2008; Wang et al. 2013). However, besides the fact that, as emphasized by Särndal et al. (1992) “Design-based inference is objective, nobody can challenge that the sample was really selected according to the given sampling design. The probability distribution associated with the design is real, not modelled or assumed”, a further advantage of a design-based approach to mapping includes obtaining consistency of maps just on the basis of the design adopted to select fairly representative or balanced samples of locations, requiring only—as a sort of nonparametric approach—mild assumptions about the population under study. In our case, IDW interpolation only requires a sort of smoothness of the surface to be interpolated with discontinuities that occur in sub-sets of measure 0 and, when dealing with finite populations of areas, regularities in the area shapes. Therefore, our IDW mapping is applicable to all the populations sharing these features. Finally, a very practical motivation is relevant in favour of a design-based approach. As stated before, in environmental and forest surveys estimation of

totals and averages avoids model-based procedures and is traditionally performed from a design-based perspective exploiting well experimented sampling schemes such as systematic grid sampling and tessellation stratified sampling (see e.g. Tomppo et al. 2010). Therefore, it would be logically inconsistent to adopt a design-based inference to perform the estimation of totals while adopting a model-based inference for mapping. That ultimately motivates the use of design-based mapping with the subsequent, practical necessity of harmonization.

Funding The authors have not disclosed any funding.

Declarations

Conflict of interest All authors declare that they have no conflict of interest.

Consent to participate All authors have consented to participate.

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1007/s00477-022-02186-2>.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Barabesi L (2003) A Monte Carlo integration approach to Horvitz–Thompson estimation in replicated environmental designs. *Metron* LXI5:355–374
- Barabesi L, Franceschi S (2011) Sampling Properties of spatial total estimators under tessellation stratified designs. *Environmetrics* 22:271–278
- Barabesi L, Franceschi S, Marcheselli M (2012) Properties of design-based estimation under stratified spatial sampling. *Ann Appl Stat* 6:210–228
- Barabesi L, Fattorini L, Marcheselli M, Pisani C, Pratelli L (2015) The estimation of diversity indexes by using stratified allocations of plots, points or transects. *Environmetrics* 26:202–215
- Bărbulescu A, Erban C, Indrean ML (2021) Computing the beta parameter in IDW interpolation by using a genetic algorithm. *Water* 13(6):863
- Basu D (1969) Role of the sufficiency and likelihood principles in sample survey theory. *Sankhya A* 31:441–454
- Chirici G, Giannetti F, McRoberts RE, Travaglini D, Pecchi M, Maselli F, Chiesi M, Corona P (2020) Wall-to-wall spatial

- prediction of growing stock volume based on Italian National Forest Inventory plots and remotely sensed data. *Int J Appl Earth Obs Geoinf* 84:101959
- Cordy CB (1993) An extension of the Horvitz–Thompson theorem to point sampling from a continuous universe. *Stat Probab Lett* 18:353–362
- Cressie N (1993) *Statistics for spatial data*. Wiley, New York
- de Gruijter JJ, ter Braak CJF (1990) Model free estimation from survey samples: a reappraisal of classical sampling theory. *Math Geol* 22:407–415
- FAO (2010) *Global Forest Resources Assessment 2010: Main Report*. FAO, Rome
- Fattorini L, Marcheselli M, Pisani C (2006) A three-phase sampling strategy for large-scale multiresource forest inventories. *J Agric Biol Environ Stat* 11:296–316
- Fattorini L, Marcheselli M, Pisani C, Pratelli L (2018a) Design-based maps for continuous spatial populations. *Biometrika* 105:419–429
- Fattorini L, Marcheselli M, Pratelli L (2018b) Design-based maps for finite populations of spatial units. *J Am Stat Assoc* 113:686–697
- Fattorini L, Marcheselli M, Pisani C, Pratelli L (2019) Design-based mapping for finite population of marked points. *Electron J Stat* 13:2121–2149
- Fattorini L, Marcheselli M, Pisani C, Pratelli L (2020) Design-based consistency of the Horvitz–Thompson estimator under spatial sampling with applications to environmental surveys. *Spat Stat* 35:100404
- Fattorini L, Marcheselli M, Pisani C, Pratelli L (2021) Design-based properties of the nearest neighbour spatial interpolator and its bootstrap mean squared error estimator. *Biometrics*. <https://doi.org/10.1111/biom.13505>
- Giraldo R, Delicado P, Mateu J (2011) Ordinary kriging for function-valued spatial data. *Environ Ecol Stat* 18:411–426
- Gong G, Mattevada S, O'Bryant SE (2014) Comparison of the accuracy of kriging and IDW interpolations in estimating groundwater arsenic concentrations in Texas. *Environ Res* 130:59–69
- Greaver TL, Clark CM, Compton JE, Vallano D et al (2016) Key ecological responses to nitrogen are related by climate change. *Nat Clim Chang* 6:836–843
- Gregoire TG (1998) Design-based and model-based inference in survey sampling: appreciating the difference. *Can J for Res* 28:1429–1447
- Gregoire TG, Valentine HT (2008) *Sampling strategies for natural resources and the environment*. Chapman and Hall, Boca Raton
- Henley S (1981) *Nonparametric geostatistics*. Applied Science Publishers, London
- Ignaccolo R, Mateu J, Giraldo R (2014) Kriging with external drift for functional data for air quality monitoring. *Stoch Env Res Risk Assess* 28:1171–1186
- Kelley HW (1990) *Keeping the land alive*, FAO Soils Bulletin 50. FAO, Rome
- Little RJ (2004) To model or not to model? Competing modes of inference for finite population sampling. *J Am Stat Assoc* 99:546–556
- Montanari GE, Cicchitelli G (2014) Sampling theory and geostatistics: a way of reconciliation. In: Mecatti F, Conti PL, Ranalli MG (eds) *Contributions to sampling statistics*. Springer, New York, pp 151–165
- Noori MJ, Hassan HH, Mustafa YT (2014) Spatial estimation of rainfall distribution and its classification in Duhok Governorate using GIS. *J Water Resource Prot* 6:75–82
- Quatemberg A (2015) *Pseudo-populations. A basic concept in statistical surveys*. Springer, Berlin
- Rao JNK, Molina I (2015) *Small area estimation*. Wiley, New York
- Särndal CE, Swensson B, Wretman J (1992) *Model assisted survey sampling*. Springer, New York
- Schreuder HT, Gregoire TG, Wood GB (1993) *Sampling methods for multiresource forest inventory*. Wiley, New York
- Smith TMF (1994) Sample surveys 1975–1990; an age of reconciliation? *Int Stat Rev* 62:5–34
- Smith TMF (2001) *Biometrika Centenary: sample surveys*. *Biometrika* 88:67–134
- Stevens DL (1997) Variable density grid-based sampling designs for continuous spatial populations. *Environmetrics* 8:167–195
- Thompson SK (2012) *Sampling*, 3rd edn. Wiley, New York
- Tobler WR (1970) A computer movie simulating urban growth in the Detroit region. *Econ Geogr* 46:234–240
- Tomppo LM, Gschwantner RE, McRoberts RE (2010) *National forest inventories: pathways for common reporting*. Springer, Heidelberg
- Wang JF, Stein A, Gao BB, Ge Y (2012) A review of spatial sampling. *Spat Stat* 2:1–14
- Wang JF, Haining R, Liu TJ, Li LF, Jiang CS (2013) Sandwich estimation for multi-unit reporting on a stratified heterogeneous surface. *Environ Plan* 45:2515–2534

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.