

## A sampling strategy for assessing habitat coverage at a broad spatial scale

Lorenzo Fattorini<sup>a,1</sup>, Marco Cervellini<sup>b,c,1</sup>, Sara Franceschi<sup>a,\*</sup>, Michele Di Musciano<sup>b,d</sup>,  
Piero Zannini<sup>b,e</sup>, Alessandro Chiarucci<sup>b,e</sup>

<sup>a</sup> Department of Economics and Statistics, University of Siena, Siena, Italy

<sup>b</sup> BIOME lab, Department of Biological, Geological and Environmental Sciences, Alma Mater Studiorum, University of Bologna, Bologna, Italy

<sup>c</sup> School of Biosciences and Veterinary Medicine, Plant Diversity and Ecosystems Management Unit, University of Camerino, Camerino, Italy

<sup>d</sup> Department of Life, Health & Environmental Sciences, University of L'Aquila, Coppito, L'Aquila, Italy

<sup>e</sup> Inter University Centre PlantDATA, Department of Biological, Geological and Environmental Sciences, Alma Mater Studiorum, University of Bologna, Bologna, Italy

### ARTICLE INFO

#### Keywords:

Design-based inference  
Italy  
Monitoring  
Precision estimation  
Two-phase sampling  
Two-stage sampling  
Simulation study

### ABSTRACT

The quantitative assessment of habitat conservation status is a major task for European Union member states in compliance with Council Directive 92/43. One goal of the European 2030 Biodiversity Strategy is the effective management of habitats that show declining trends. While various approaches have been adopted for national assessments, there is no consensus on how to achieve common statistically sound estimates of the criteria indicated by the EU Directive for the evaluation of the status and trend of habitat types. Here, we present an adaptive monitoring approach based on a two-phase sampling scheme to estimate the coverage of EU terrestrial habitat types, which is one of the four criteria indicated by the Habitats Directive. We used 9 habitats distributed among different EU member states choosing Italy as a case study. The development of the methodological approach is described, and a simulation study was performed to check the precision of the coverage estimators accounting for the lack of sampled data (nonresponse treatment), subregions and sustainable sampling effort. We found that our two-phase sampling approach has the potential to increase precision in estimating the coverage of habitat types (approximated at 1 ha cell size) with respect to the precision achieved by simple random sampling without replacement, which is the simplest sampling approach. Adopting a small sampling fraction ( $\leq 0.04\%$ ) of the survey area, the relative standard errors ranged from 7 to 15% for common habitats whose presence is strongly correlated with the habitat suitability scores furnished by an expert team. In the challenging context of a “mandated” monitoring type, our approach provides sound statistical estimates of habitat coverage with the possibility of applying a standardised and transferable sampling scheme that is easily repeatable over time.

### 1. Introduction

The decline of habitat structure and functioning is becoming increasingly relevant worldwide, not only because of the consequent dramatic biodiversity loss but also due to the quantitative and qualitative decrease in various ecosystem services (Martinez-Harms et al., 2015; Mulder et al., 2015; Keyes et al., 2021). The need to know how habitats are changing is becoming a pivotal challenge in applied ecology, and science-based information provides solid strategic information for global decision-making. In this direction, long-term ecological research is necessary to fulfil ecological and social goals, and it should be founded on well-designed long-term adaptive (sensu Lindenmayer and

Likens, 2009, 2018; Lindenmayer et al., 2020) monitoring approaches (Smith and Gray, 2021; Cowles et al., 2021).

In Europe, the conservation of EU habitat types within and outside of the world's largest coordinated network of protected areas (Natura 2000 Network, hereafter N2K; European Commission, 2021) is one of the main targets of EU Directive 92/43. Furthermore, the EU Biodiversity Strategy 2030 targets the enforcement of the Habitats Directive (HD) by enlarging the N2K and by improving the effective management of sites and habitats that show declining trends (European Commission, 2020). Under this framework, it is essential to achieve quantitative and affordable measures of the status and trends concerning habitat area and quality.

\* Corresponding author.

E-mail addresses: [lorenzo.fattorini@unisi.it](mailto:lorenzo.fattorini@unisi.it) (L. Fattorini), [marcocervellini@gmail.com](mailto:marcocervellini@gmail.com) (M. Cervellini), [sara.franceschi@unisi.it](mailto:sara.franceschi@unisi.it) (S. Franceschi), [michele.dimusciano@univaq.it](mailto:michele.dimusciano@univaq.it) (M. Di Musciano), [piero.zannini2@unibo.it](mailto:piero.zannini2@unibo.it) (P. Zannini), [alessandro.chiarucci@unibo.it](mailto:alessandro.chiarucci@unibo.it) (A. Chiarucci).

<sup>1</sup> Equally contributing authors.

<https://doi.org/10.1016/j.ecolind.2022.109352>

Received 20 February 2022; Received in revised form 11 August 2022; Accepted 21 August 2022

1470-160X/© 2022 The Authors. Published by Elsevier Ltd. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

The HD describes a “habitat” as an area of a species’ occurrence defined by geographic, abiotic, and biotic features. The ecological literature provides a variety of alternative habitat definitions contributing to a detailed qualitative description of the autecological species-related habitat concept, while an operational and standardized quantitative characterization of this concept is still lacking (e.g., SYapp, 1922; Hall et al., 1997; Davies et al., 2004; Mitchell, 2005; Drakou et al., 2011; Cervellini et al., 2021). In this context, Fahrig, 2013 recently proposed the “habitat patch concept”. This concept assumes that habitat patch boundaries contain and delimit biological populations and communities and raises the problem of how to delineate and measure ecologically relevant habitat patches. Although new approaches to measure and model the terrestrial area of a habitat are discussed (e.g., area of occupancy “AOO” and area of habitat “AOH”, see Álvarez-Martínez et al., 2018; Brooks et al., 2019), complete habitat type mapping is still a challenge, particularly when the habitat classification and the related conservation goals differ among organizations or regulations (e.g., IUCN and HD). In summary, the debate between categorical and dynamic habitat maps is ongoing (discrete classes of habitat vs. continuous values, see Coops and Wulder, 2019), and a univocal and operational definition of the parameter “area of habitat” is still lacking.

Formally, the evaluation of the conservation status of EU habitat types is based on four criteria: “range”, “area”, “habitat quality” and “pressures and threats”, but a standardized approach for ecological monitoring at the EU scale is still lacking (Ellwanger et al., 2018; Lengyel et al., 2018; Delbosc et al., 2021). This gap is quite understandable in the complex context of habitat recognition and mapping. Despite methodological problems, monitoring the conservation status of habitats listed in Annex I is mandatory according to Article 11 of EU Directive 92/43. Each member state (MS) must submit a national report to the European Commission every six years on the implemented measures and their effectiveness (Art. 17 HD) based on monitoring results. Most EU countries are producing six-year reports based on expert-based assessments or supposed complete censuses. For instance, in Italy, during the past four reporting cycles (1994–2018), the Institute for Environmental Protection and Research (ISPRA) provided the European Commission with a habitat conservation status assessment for both national and biogeographical regions by merging the data independently gathered by the 21 Italian regions and autonomous provinces. The habitat monitoring actions performed by these local public agencies or institutions were based on standardized guidelines concerning “how” to survey in the field (Angelini et al., 2016) but without indications about “where” (i.e., sampling scheme) and “how much” to survey (i.e., sampling effort). These omissions made it extremely difficult to fully merge the data and perform statistical inferences on countrywide habitat population and each biogeographical region and to quantify and detect changes or trends in the targeted criteria between the different reporting cycles. Ellwanger et al., 2018 showed that none of the surveyed MSS made a theoretical statement on the statistical strength of the adopted monitoring approaches, highlighting the need for better sampling and assessment approaches.

Developing an adaptive and statistically sound long-term monitoring plan is becoming pivotal to establish how data should be collected and to produce standardized, reliable estimates for given parameters (e.g., area). Data from such a sample survey should be (a) representative of the population under investigation and (b) information-rich to reduce uncertainty about inferences (Foster, 2020). Achieving these outcomes at the continental scale becomes particularly complex in a “mandated” monitoring plan (Lindenmayer and Likens, 2010) as that imposed by the HD. In this context, it is crucial to provide a standardized sampling strategy (Delbosc et al., 2021) to guarantee the following three properties are considered generally relevant in determining the scientific quality of biodiversity monitoring (Lengyel et al., 2018): (i) a sound and feasible sampling scheme, (ii) a good trade-off between sampling effort and the precision of the resulting estimators, and (iii) appropriate statistical analysis to detect changes or trends.

Here, we develop a two-phase sampling strategy to estimate quantities approximating the area of terrestrial habitats, which is one of the four criteria indicated by the HD for evaluating the conservation status of EU habitat types. After conceptual analysis and methodological development of the strategy, a simulation study was performed to check and compare the precision of the proposed estimators for nine selected habitats, seven of which were distributed among different EU countries (EIONET, 2022) under several potential and real constraints (e.g., critical aspect emerged during the previous four reporting cycles) that may arise during the surveys (e.g., presence of auxiliary information, non-responses). This approach was developed for the territory of a single country, namely, Italy, but can be rescaled to the territory of any other country or to the entire European Union.

## 2. Materials and Methods

### 2.1. Study region

The study region was the surface area within the administrative borders of the Italian state. It spans 301,328.46km<sup>2</sup> and is covered by 3,491 quadrats of 10km × 10km of the grid used for reporting HD data by European member states (European Environment Agency, 2013; Cervellini et al., 2020). Therefore, the total area of the grid overlapping the country is 349,100km<sup>2</sup> and includes some parts outside Italian borders and the sea (see Fig. 1). The presence of 124 habitats in the study region was stated in the ISPRA report (ex art. 17 HD), together with the number of quadrats (hereafter denoted *M*) in which they were present (see Table 1). We used this information as the starting point to construct the sampling strategy for estimating quantities approximating the area of each habitat in the study region.

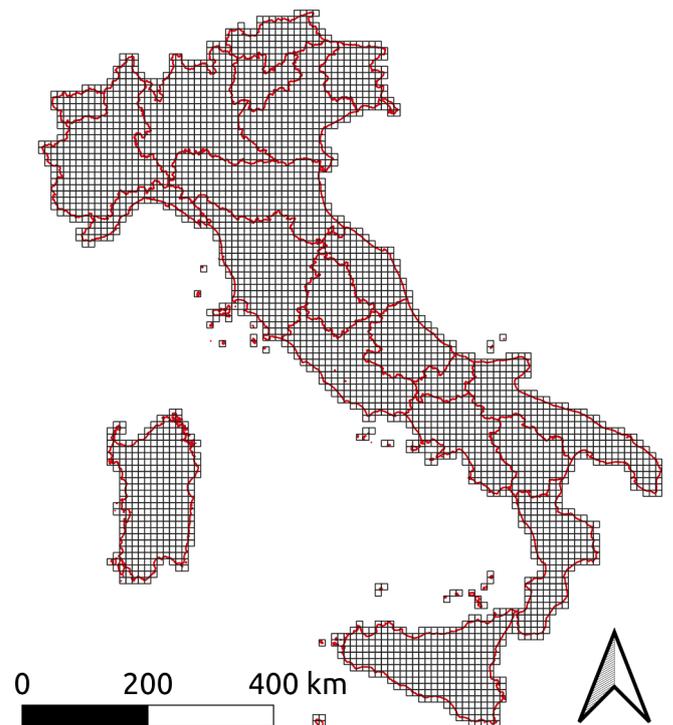


Fig. 1. The study region is covered by 3,491 quadrats of 10km × 10km (black border and white filling), some of which lie partially outside Italian borders. The administrative borders of the Italian regions are represented by the red line. The map was obtained by intersecting the standard reference grid of quadrats 10km × 10km provided by the European Environmental Agency (European Environment Agency, 2013) for the Italian surface area with the administrative borders of the Italian regions (ISTAT, 2022).

**Table 1**

Number of 10km × 10km quadrats (*M*) of habitat presence for the 124 habitats present in the study region listed by their HD codes. Habitats highlighted in bold are adopted in the simulation study (see also Fig. 4)

code	<i>M</i>	code	<i>M</i>	code	<i>M</i>	code	<i>M</i>
1150	217	3250	392	6410	342	9180	514
1210	552	3260	563	6420	297	9190	18
1240	352	3270	618	6430	937	91AA	1267
1310	222	3280	355	6510	1031	91B0	55
1320	23	3290	221	6520	427	91D0	60
1340	3	4030	303	7110	98	91	1091
1410	247	4060	556	7120	2	91F0	325
1420	199	<b>4070</b>	<b>262</b>	7140	246	91H0	123
1430	169	4080	236	7150	80	91K0	188
1510	59	4090	140	7210	122	91L0	448
2110	359	<b>5110</b>	<b>48</b>	7220	307	91M0	619
2120	260	5130	499	7230	318	<b>9210</b>	<b>516</b>
2130	32	5210	277	7240	62	9220	141
2160	6	5220	16	8110	362	9250	33
2210	193	5230	82	8120	403	9260	1118
2230	322	5310	6	8130	481	<b>92A0</b>	<b>1387</b>
2240	173	5320	193	8210	1134	92C0	42
2250	212	5330	930	8220	551	92D0	505
2260	133	5410	20	8230	262	9320	247
2270	188	5420	32	8240	160	<b>9330</b>	<b>390</b>
2330	8	5430	101	8310	757	9340	318
3110	8	6110	436	8320	43	9350	4
3120	74	6130	90	8330	153	9380	41
3130	615	6150	386	8340	128	<b>9410</b>	<b>368</b>
3140	326	6170	552	9110	402	<b>9420</b>	<b>442</b>
3150	842	6210	1473	<b>9120</b>	3	9430	52
3160	39	<b>6220</b>	<b>1566</b>	9130	334	9510	28
3170	301	6230	548	9140	70	9530	103
3220	394	6240	45	9150	143	9540	195
3230	87	62A0	225	9160	209	9560	6
3240	597	6310	159	9170	1	9580	20

**4070** - Bushes with *Pinus mugo* and *Rhododendron hirsutum* (*Mugo-Rhododendretum hirsuti*); **5110** - Stable xero-thermophilous formations with *Buxus sempervirens* on rock slopes (*Berberidion pp*); **6220** - Pseudo-steppe with grasses and annuals of the *Thero-Brachypodietea*; **9120** - Atlantic acidophilous beech forests with *Ilex* and sometimes *Taxus* in the shrublayer (*Quercion robori-petraeae* or *Ilici-Fagenion*); **9210** - Apennine beech forests with *Taxus* and *Ilex*; **92A0** - *Salix alba* and *Populus alba* galleries; **9330** - *Quercus suber* forests; **9410** - Acidophilous *Picea* forests of the montane to alpine levels (*Vaccinio-Piceetea*); **9420** - Alpine *Larix decidua* and/or *Pinus cembra* forests. The entire set of habitat codes along with the related number (*M*) of 10km × 10km quadrats of habitat presence was extracted from the distribution habitat maps provided by the official Eionet European Central Data Repository (CDR, 2022) for the progress reports and implementation of Article 17(HD).

## 2.2. Survey arrangement

It was difficult to accurately delineate the ground distribution of most habitats. While recording the size of patches containing the habitat was challenging, recording the presence of the habitat in fixed-area units of adequate size was straightforward. Therefore, we partitioned the *M* quadrats in which the habitat was present into a grid of  $K = 100m \times 100m$  square cells (1ha). This cell size was considered a good compromise between the need to perform a complete ecological and cost-effective habitat survey within the cell and an appropriate approximation of the total area. The cells completely outside the Italian borders or completely overlapping with the sea were discarded; the remaining cells constituted the initial population  $U_0$ . Moreover, to avoid surveying cells where the habitat presence was impossible, we needed to quantify the chance of habitat presence in the cells. Therefore, for each cell  $j \in U_0$ , the ISPRA Group for Terrestrial Habitat Monitoring and Conservation calculated a value  $x_j \geq 0$ , referred to as the habitat suitability score (HSS). Scores equal to 0 were assigned to cells where the habitat presence was impossible, and these cells were discarded (see Section 2.3). Then, the target population to be surveyed was constituted by the set  $U \subset U_0$  of *N*

cells in which  $x_j > 0$ , i.e., habitat presence was considered possible (see Fig. 2).

The target parameter under estimation was established to be the number of cells *Y* in which the habitat is present. In practice, *Y* constitutes the total area of the cells that cover the habitat at a grain size of 1ha, and as such, it is referred to as habitat coverage. To estimate *Y*, we introduced a survey variable indexing the presence/absence of the habitat in the cells, i.e., for each cell  $j \in U$ ,  $y_j$  was set to be 1 if the habitat was present in the cell and 0 otherwise. In this way, the target quantity *Y* was the population total, i.e.,

$$Y = \sum_{j \in U} y_j.$$

We aimed to select samples with cells evenly spread throughout the region of habitat presence, thus achieving spatial balance (e.g., Grafström and Lundström, 2013; Brown et al., 2015). Moreover, to maximize the likelihood of encountering the habitat within the selected cells, HSS values were adopted as auxiliary information to guide cell sampling. Owing to the large effort that may be required for detecting the habitat presence with 1ha cells, we established a maximum sampling fraction of 0.04%.

A further source of auxiliary information was the habitat presence in some cells from recent investigations. In particular, we used recently available habitat maps representing spatial polygons with previously validated habitat presence (see the methodology for producing the habitat maps in Carli et al., 2020) and the spatial polygons representing all the Italian protected sites within the N2K network (MiTE, 2021). The N2K is the largest coordinated network of protected areas in the world, extending across all the 28 EU countries and designated under the HD. Specifically, for each cell  $j \in U$ ,  $h_j$  was set to 1 if it was known that the habitat was present in the cell and 0 otherwise. Accordingly, the population total of this variable was

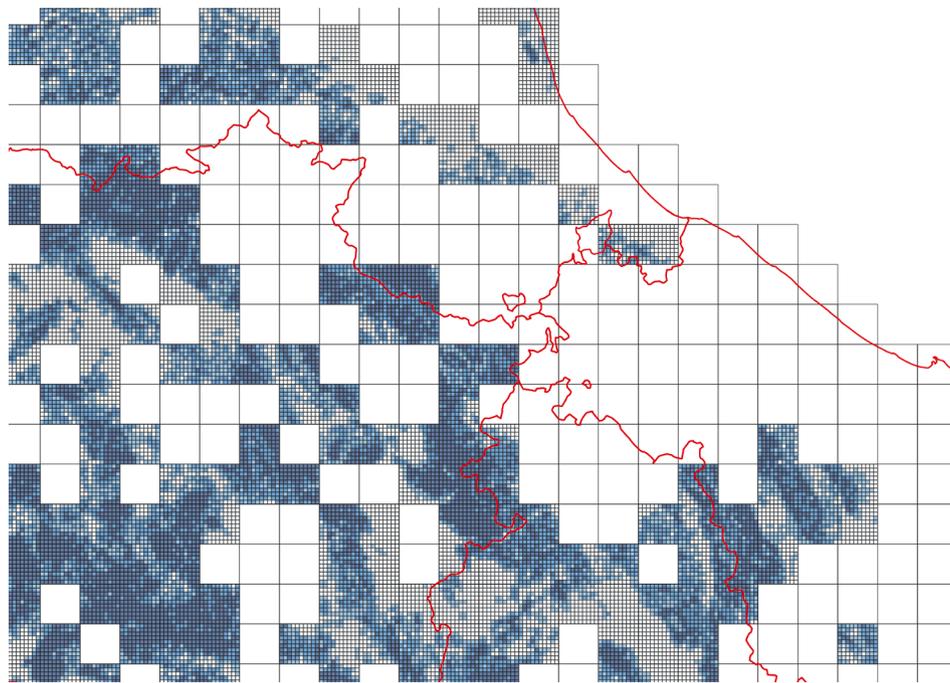
$$H = \sum_{j \in U} h_j$$

and constituted the extent of the region in which habitat presence was certain and was the lower bound for any estimate of *Y*.

## 2.3. Determination of habitat suitability scores (HSSs)

HSSs were determined based on a set of variables correlated with the presence/absence of the habitat in the cells (i.e., survey variable) and available for all cells in the study region. Thus, from variables for the habitat characterization provided by the official manual for habitat monitoring (Angelini et al., 2016) and available (informatic layers) at the national scale, ISPRA experts first selected the following (Table 2): (i) land use types (CLC, 2018), (ii) exposure (derived from DEM20), (iii) altitude (derived from DEM20 - SINAnet 2020), (iv) slope (derived from DEM20), (v) hydrographic network (HydRet - SINAnet 2020) and (vi) distance to the coastline (ISTAT, 2020). Each predictor included a set of classes derived from the structure of the data (e.g., CLC, 2018) or established based on expert knowledge at the national scale (e.g., distance to the coastline for coastal habitats). For each habitat type, experts assigned one of the following weights to each predictor class: "0" (null suitability for habitat presence), "0.5" (intermediate suitability for habitat presence), or "1" (high suitability for habitat presence).

For each cell  $j \in U_0$ , we then obtained an HSS score  $0 \leq x_j \leq 1$  multiplying all the weights associated with each predictor such that a score resulted in 0, i.e., no possibility of habitat presence, if at least one of the weights was 0. All the predictors listed in Table 2 were used to ecologically characterize the habitat types adopted for this study. Notably, if the predictor was considered ecologically irrelevant to define the distribution of a specific habitat type, we assigned a value of 1 to all categories of the "irrelevant variable", thus not affecting the final suitability score (as the suitability score is the result of multiplication, and multiplying by 1 does not change the suitability score). To improve



**Fig. 2.** An example of the population of cells with HSSs score > 0 to be sampled. The colour scale from light blue to blue indicates the increasing value of the HSS scores. White-coloured cells had an HSS equal to 0 and were discarded. The administrative borders of the Italian regions are represented by the red line.

**Table 2**  
Summary of the environmental predictors selected for determining habitat suitability scores (HSSs). The reported on-line data sources were visited in July 2020.

Description	Type	Data source
Corine Land Cover	Raster	CLC 2018. Version is v.2020_20u1. <a href="https://land.copernicus.eu/pan-european/corine-land-cover/clc2018?tab=download">https://land.copernicus.eu/pan-european/corine-land-cover/clc2018?tab=download</a>
Digital Elevation Model	Raster	DEM20. Digital Elevation Model. Rete del Sistema Informativo Nazionale Ambientale. SINAnet. <a href="http://www.sinanet.isprambiente.it/it/sia-ispra/download-mais/dem20/view">http://www.sinanet.isprambiente.it/it/sia-ispra/download-mais/dem20/view</a>
Hydrographic network	Shape file	HydRet. Hydrographic network. Rete del Sistema Informativo Nazionale Ambientale. SINAnet. <a href="http://www.sinanet.isprambiente.it/it/sia-ispra/download-mais/reticolo-idrografico/view">http://www.sinanet.isprambiente.it/it/sia-ispra/download-mais/reticolo-idrografico/view</a>
Coastal Line	Shape file	ISTAT 2020-Sezioni di Censimento Litoranee. Linea litoranea. <a href="https://www.istat.it/it/archivio/137341">https://www.istat.it/it/archivio/137341</a>

the chance of sampling cells in which the habitat was present we increased the HSS of cells falling within a known habitat coverage polygon (Carli et al., 2020) by multiplying it by 8. Moreover, to further improve the chance of sampling cells within the N2K network polygons, their HSS value was multiplied by 1.25. Consequently, the HSS of a cell located both within a habitat coverage polygon and a N2K polygon that varied originally from 0 to 1 was multiplied by 8 and then by 1.25. After these rescaling operations, HSSs were in the range 0–10.

#### 2.4. Sampling and estimation

Cell sampling was performed separately for each habitat. We adopted a two-phase scheme, with the first phase comprising two stages (Sections 1–4 of the Supporting Information file). The complexity of the scheme was due to the necessity of a final sample of cells evenly distributed among and within the quadrats and having high HSS values.

The first stage of the first phase was performed to evenly distribute

selected quadrats throughout the study region. For this purpose, the set of the  $M$  quadrats in which the habitat was present was partitioned into  $m$  clusters of neighbouring quadrats, referred to as the q-blocks. The number of quadrats per q-block was established to ensure that blocks had approximately the same number of quadrats. Partition was performed adopting the k-means algorithm (k-means clustering - “stats” package R Core Team, 2020). This algorithm was originally proposed by Hartigan and Wong, 1979. For this purpose, the algorithm needed the number of clusters  $m$ , the spatial coordinates of the quadrat centroids and the maximum number of iterations allowed, which was established to be 100,000. Subsequently, in accordance with the sampling scheme referred to as the one-per-stratum sampling (e.g., Breidt, 1995), henceforth OPSS, a quadrat was selected in each block with probabilities proportional to the HSS totals within blocks. Because a unique quadrat was selected within clusters of contiguous quadrats, OPSS ensured that samples of quadrats were evenly spread throughout the study region. Moreover, because the selection was performed with probabilities increasing with the HSS totals within quadrats, the quadrats with high HSSs had a greater chance of being selected.

The number  $m$  of quadrats to select from the  $M$  was established by the following function of  $M$

$$m = \begin{cases} M & \text{if } M \leq 10 \\ [9.0756303 + 0.0924369M] + 1 & \text{if } 10 < M < 1200 \\ [0.1M] & \text{if } M \geq 1200 \end{cases} \quad (1)$$

where  $[x]$  is the integer part of  $x$ . The algorithm was used to adjust the sampling effort with respect to the total number of quadrats  $M$ , avoiding excessive effort when the number of quadrats was large. In particular, the algorithm established that no selection was performed if the number of quadrats was smaller than 10, in which case all quadrats were included in the sample, while the percentage of selected quadrats decreased linearly from 100% when  $M = 10$  to 10% when  $M = 1200$ , remaining equal to 10% for any  $M$  greater than 1200.

The second stage of the first phase was performed to evenly spread the selected cells within the quadrats selected in the first stage. For this purpose, OPSS was once again performed within the selected quadrats. Because the cells were arranged in a regular grid of size  $100m \times 100m$ ,

there was no need for time consuming clustering algorithms to determine clusters of neighbouring cells. The  $m$  quadrats selected in the first stage were partitioned into  $k = 25$  quadrat blocks of  $20 \times 20$  cells, referred to as c-blocks (see Fig. 3). The blocks where the habitat was absent were discarded, and a cell was randomly selected within each of the remaining c-blocks in accordance with the OPSS scheme with probabilities proportional to the HSSs to ensure that cells with high HSSs had a higher chance of being selected. Therefore, a maximum of 25 cells was selected within each selected quadrat at the end of the first phase.

Finally, the second phase was performed to reduce the sampling effort from a maximum of  $k = 25$  cells per quadrat to a maximum of  $\bar{n} = 4$  cells. Because the spatial balance of selected cells within quadrats was already achieved by the use of OPSS in the second stage of the first phase, the spatial component was ignored and only the HSSs were considered to ensure that cells with high HSSs had a higher chance of being selected. Accordingly, a sample of  $\bar{n} = 4$  cells was selected with probability proportional to the HSS of the cells selected in the first phase within each quadrat. Selection was performed with the Sampford algorithm (Sampford, 1967). If the cells selected in a quadrat at the end of the first phase were less than or equal to 4, second-phase selection was not carried out, and all the cells were included in the final sample.

Once the final sample was achieved,  $Y$  was estimated by means of the double expansion (DE) estimator (e.g., Särndal et al., 1992, Section 9.3)  $\hat{Y}_{(2)DE}$  given by equation (SM.9). The DE estimator was adopted because it is able to handle the complexities involved in using multi-phase sampling schemes. The DE estimator was design-unbiased with the design-based variance given by equation (SM.11). If the habitat presence was known for some cells in the population, we exploited this additional information by means of the difference (DIF) estimator (e.g., Särndal et al., 1992, Section 6.3)  $\hat{Y}_{(2)DIF}$  given by equation (SM.13). The DIF estimator was simply a modification of the DE estimator to include the additional information in the estimation criterion. It has been recently used in biodiversity surveys to improve species richness estimation (Chiarucci et al., 2018). The DIF estimator was design-unbiased with design-based variance given by equation SM.14. In our case, the DIF estimator had the appealing property to providing consistent results in that the estimate was never smaller than the number of cells  $H$  in which the habitat presence was known, which obviously should constitute a

lower bound for any estimator. This feature was not ensured by the DE estimator.

Notably, auxiliary information from which to construct HSSs and previous knowledge of the cells with habitat presence are not essential for executing the strategy. If no auxiliary information is available, all the HSSs are set to be equal such that all the quadrats have an equal chance to be selected, and only the even spread of the selected quadrats and cells is ensured by the OPSS. Moreover, if no previous knowledge of habitat presence is available, all the  $h_j$  are set equal to 0 such that DIF and DE estimators coincide.

The variances of the DE estimator  $\hat{Y}_{(2)DE}$  and the DIF estimator  $\hat{Y}_{(2)DIF}$ , were estimated using the Hansen–Hurwitz (HH)-like variance estimators (e.g., Wolter, 2007)  $V_{DE}^2$  and  $V_{DIF}^2$  according to equations (SM.12) and (SM.15), respectively.

## 2.5. Nonresponse treatment

For cells that are impossible to reach in the field, due to topographic constraints, it was impossible to verify the habitat presence. We thus treated these missing values as nonresponses, and we followed the design-based suggestion by Fattorini et al., 2013, i.e., we corrected the DE and DIF estimators performed on the respondent sample by nonresponse calibration weighting (Haziza et al., 2010). The purpose was to increase the estimates achieved from the respondent sample to reduce the downwards bias invariably induced by nonresponses. By this approach, nonresponses were viewed as fixed characteristics of the cells, without attempting any model to explain them. This approach was performed by introducing for each cell a response indicator  $z_j$  that was equal to 1 if it was possible to reach and explore and 0 otherwise (see Section 7 in the Supporting Information file for information on nonresponse treatment in sample surveys).

If the DE criterion was adopted, we then calibrated the estimator computed on the respondent sample by the estimator  $\hat{Y}_{(2)DE-CAL}$  according to equation (SM.19). If the DIF criterion was adopted, calibration was performed by the estimator  $\hat{Y}_{(2)DIF-CAL}$  in equation (SM.23). Section 7 in the Supporting Information file provided the conditions under which the two calibrated estimators reduced the downwards bias and turned out to be approximately unbiased. The condition was that the relationships between the  $y_j$ s and  $x_j$ s in the case of the estimator  $\hat{Y}_{(2)DE-CAL}$  or that between the  $d_j$ s and  $x_j$ s in the case of  $\hat{Y}_{(2)DIF-CAL}$  were similar in respondent and nonrespondent cells (see also Fattorini et al., 2013).

The variances of  $\hat{Y}_{(2)DE-CAL}$  and of  $\hat{Y}_{(2)DIF-CAL}$  were estimated by the HH-like variance estimators  $V_{DE-CAL}^2$  and  $V_{DIF-CAL}^2$  according to equations (SM.22) and (SM.24), respectively.

## 2.6. Coverage estimation within subregions

Coverage estimation was performed for the entire Italian surface area and for the three biogeographical regions partitioning the study region (Cervellini et al., 2020) as well as for the portion of cells located within the Natura 2000 Network. For this purpose, we introduced for each cell  $j \in U$  an indicator  $u_{g,j}$  that was equal to 1 if the cell was in the subregion  $g$  of interest and 0 otherwise. Then, we adopted the estimator  $\hat{Y}_{(2)DE}$  or  $\hat{Y}_{(2)DIF}$  in the case of complete samples or the estimator  $\hat{Y}_{(2)DE-CAL}$  or  $\hat{Y}_{(2)DIF-CAL}$  in the case of nonresponses together with their corresponding variance estimators simply by multiplying the  $y_j$ s by  $u_{g,j}$ s. This strategy has been widely adopted in sample surveys when estimating totals in particular sectors of populations and is usually referred to as domain estimation (see Särndal et al., 1992, Chapter 10 and Section 8 in the Supplementary Information file for more details).

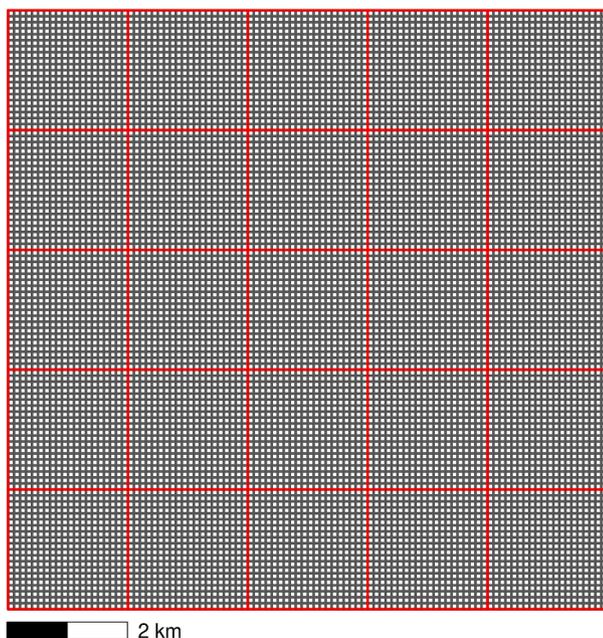


Fig. 3. First stage selected quadrat partitioned into  $k = 25$  square c-blocks (c-block edges in red) of 400 cells (cell edges in black).

2.7. Simulation study

A simulation study was performed to check and compare the precision of the estimation strategies corresponding to the estimators  $\hat{Y}_{(2)DE}$  and  $\hat{Y}_{(2)DIF}$ . To be realistic, we selected nine habitats from those listed in Table 1 for which HSS values were available. The selected habitats showed different spatial distributions, representing the three levels of spatial occurrences within quadrats (i.e.,  $M \leq 10, 10 < M < 1200$  and  $M \geq 1200$ , see Algorithm 1), with the very rare habitat 9120, the scattered habitat 5110, the more common habitats 4070, 9210, 9330, 9410 and 9420 and the very common habitats 6220 and 92A0 (see Fig. 4).

For these habitats, we standardized HSSs in the range 0–1 by dividing each HSS by 10. The number of cells  $N$  with HSS values greater than 0 indicated the territory (in *ha*) where the habitat presence was possible, i.e., the extent of the survey area. Based on these territories and from previous information on habitat validated presence (Carli et al., 2020), we attempted to provide realistic coverages for each habitat, establishing the number of cells  $Y$  where the habitat was present. The resulting coverages ranged from a minimum of 0.4% for habitat 9120 to a maximum of 7.9% for habitat 9410 (see Table 3).

We then generated habitat presence by sorting the cells with respect to their HSSs from the greatest to the smallest value and assigning  $y_j = 1$  to the first  $Y$  cells and  $y_j = 0$  to the remaining  $N - Y$ . In practice, we generated presence in cells with the greatest HSSs. Moreover, knowledge of habitat presence from past investigations was incorporated by assigning  $h_j = 1$  to cells for which  $y_j = 1$ , which were also present in the ISPRA polygons, and assigning  $h_j = 0$  to the remaining cells. Once the  $y_j$ s and  $h_j$ s were generated, their totals  $Y$  and  $H$  were determined. All these values remained fixed throughout the simulation runs as fixed characteristics because in design-based approaches, uncertainty stems only from sampling.

For each habitat, we independently performed  $R = 100,000$  two-phase selections of cells following the sampling scheme described in subSection 2.4. In the second phase, we selected a maximum number of  $\bar{n} = 1, 2, 3, 4$  cells. Because the final sample size  $n$  was a random variable depending on the number of cells with positive HSSs within the selected c-blocks, the expected sample size (ESS) was empirically computed as

$$ESS = \frac{1}{R} \sum_{r=1}^R n_r$$

where  $n_r$  is the size of the final sample selected at the  $r$ -th simulation run. Moreover, the expected fraction of habitat presence in the sample (EPS) was empirically computed as

$$EPS = \frac{1}{R} \sum_{r=1}^R \frac{H_r}{n_r}$$

where  $H_r$  is the number of cells with habitat presence in the final sample. The EPS values were compared with those expected under simple random sampling without replacements (SRSWOR) that coincided with the fraction of cells with habitat presence in the population, i.e.,  $p = Y/N$ .

Moreover, for each sample, the estimators  $\hat{Y}_{(2)DE}$  and  $\hat{Y}_{(2)DIF}$  were computed from the sample data together with their variance estimators  $V_{DE}^2$  and  $V_{DIF}^2$ . At the end of the procedure, for each habitat and both estimators, we determined  $R$  coverage estimates,  $\hat{Y}_1, \dots, \hat{Y}_R$ , and the corresponding variance estimates, i.e.,  $V_1^2, \dots, V_R^2$ , from which we derived the relative standard error estimates  $RSE_1, \dots, RSE_R$  with  $RSE_r = V_r / \hat{Y}_r$  for  $r = 1, \dots, R$ . Finally, the confidence interval at the nominal level of 0.95 was achieved by  $\hat{Y}_r \pm 2V_r$ . For each habitat and each estimator, the two collections constituted the Monte Carlo distributions of the abundance estimator and of its relative standard error estimators. The collections mimicked the unknown corresponding distributions and were adopted to empirically determine the theoretical properties.

Accordingly, from the resulting Monte Carlo distributions achieved by simulation, the expectation and the variances of the abundance estimator were empirically determined as follows:

$$E = \frac{1}{R} \sum_{r=1}^R \hat{Y}_r$$

and

$$Var = \frac{1}{R} \sum_{r=1}^R V_r^2.$$

From these quantities, the relative bias  $RB = (E - Y)/Y$  and the relative standard error  $RSE = \sqrt{Var}/Y$  were determined. We then tested the design effect (e.g., Särndal et al., 1992, Section 2.10) in terms of the RSE. In practice, the RSEs of the two estimators were compared with those achieved by the Horvitz-Thompson (HT) estimator under SRSWOR with sample size ESS, i.e.,

$$RSE_{SRSWOR} = \sqrt{\frac{N - ESS}{N \times ESS} \frac{1 - p}{p}}.$$

Moreover, the expectation of the relative standard error estimators was achieved as follows:

$$ERSEE = \frac{1}{R} \sum_{r=1}^R RSE_r.$$

Finally, the actual coverage of the nominal 0.95 confidence intervals C95 was obtained as the fraction of the intervals containing the true coverage  $Y$ . Additionally, both estimators ensured design-unbiasedness. Therefore, their RB values were theoretically known to be 0, and their empirical counterpart RBs were considered only to confirm the reliability of the simulation study.

To check the effectiveness of bias reduction in the presence of non-responses, a further simulation study was performed. Nonresponses were artificially generated from the populations adopted in the previous study assigning a dichotomous index  $r_j = 1$  if it was possible to reach the cell  $j$  and  $r_j = 0$  otherwise. Nonresponses were considered impossible, i.e.,  $r_j = 1$  for those cells such that  $h_j = 1$ , i.e., cells where habitat presence was known from previous investigations, which obviously implied

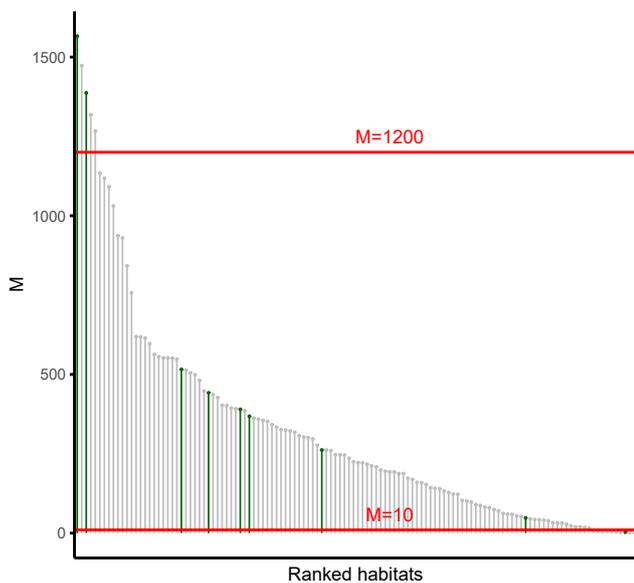


Fig. 4. Bar graph showing the ranked  $M$  for the 124 habitats in the study region reported in Table 1. Horizontal lines denote the three levels of occurrence ( $M \leq 10, 10 < M < 1200$  and  $M \geq 1200$ ). Habitats highlighted in green were adopted in the simulation study.

**Table 3**

Monte-Carlo performance of double expansion and difference estimators of coverage compared with the Horvitz-Thompson estimator under simple random sampling and performance of relative standard error estimators for nine habitats in the study region

Habitat	$\bar{n}$	ESS	EPS (%)	RSE (%)			ERSEE (%)		C95 (%)	
				SRSWOR	DE	DIF	DE	DIF	DE	DIF
<b>4070</b> survey area=1,948,964ha coverage=110,000 ha (5.6%) known coverage=72,559ha (66.0%) ss mean=0.10 corr. presence/ss=0.88 M = 262, m = 34	1	34.0	33.5	70.1	31.4	28.8	38.7	21.1	94.0	73.5
	2	68.0	33.5	49.6	23.4	20.6	28.3	17.9	94.9	81.4
	3	102.0	33.0	40.5	20.2	17.1	23.5	15.6	95.1	85.0
	4	136.0	32.3	35.1	18.4	15.0	20.5	13.9	94.9	86.8
<b>5110</b> survey area=446,098ha coverage=4,000ha (0.9%) known coverage=1,225ha (30.6%) ss mean=0.05 corr. presence/ss=0.84 M = 48, m = 14	1	14.0	7.1	281.0	165.6	163.8	66.2	32.3	65.3	43.5
	2	28.0	7.1	198.7	119.5	117.0	71.3	38.8	75.2	60.7
	3	42.0	7.1	162.2	95.8	92.9	67.4	40.3	77.9	66.9
	4	56.0	6.9	140.5	83.7	80.6	62.3	39.9	82.7	70.3
<b>6220</b> survey area=12,820,526ha coverage=350,000ha (2.7%) known coverage=87,465ha (25.0%) ss mean=0.03 corr. presence/ss=0.89 M = 1566, m = 157	1	157.0	33.7	47.6	14.5	14.0	17.5	15.6	97.0	95.7
	2	313.9	33.6	33.7	10.8	10.2	12.4	11.1	97.1	96.1
	3	470.4	33.2	27.5	9.3	8.6	10.2	9.1	96.5	95.9
	4	626.6	32.3	23.9	8.4	7.7	8.8	7.9	96.0	95.5
<b>9120</b> survey area=28,993ha coverage=125ha (0.4%) known coverage=121ha (98.8%) ss mean=0.08 corr. presence/ss=0.84 M = 3, m = 3	1	3.0	17.7	877.4	122.6	43.0	43.1	0.4	46.3	0.5
	2	6.0	17.6	620.4	85.9	32.3	63.1	0.9	73.6	1.2
	3	9.0	17.6	506.5	67.2	24.8	69.3	1.0	89.1	1.6
	4	12.0	15.2	438.6	59.5	21.8	66.7	1.3	90.6	2.3
<b>9210</b> survey area=5,058,912ha coverage=240,000ha (4.7%) known coverage=221,279ha (92.2%) ss mean=0.12 corr. presence/ss=0.97 M = 516, m = 57	1	57.0	31.1	59.4	20.8	13.7	28.7	5.9	97.9	42.5
	2	114.0	31.2	42.0	16.1	9.8	20.5	5.8	97.8	57.7
	3	171.0	31.1	34.3	14.1	8.0	16.8	5.4	97.1	65.0
	4	228.0	30.5	29.7	13.1	7.0	14.6	5.1	96.4	69.2
<b>92A0</b> survey area=12,683,736ha coverage=130,000ha (1.0%) known coverage=77,418ha (59.6%) ss mean=0.07 corr. presence/ss=0.80 M = 1387, m = 139	1	139.0	26.2	83.3	17.3	13.7	18.7	13.2	95.5	91.1
	2	278.0	26.3	58.9	13.3	10.0	13.3	9.6	94.7	92.6
	3	417.0	25.9	48.1	11.7	8.5	10.9	7.9	93.2	92.6
	4	556.0	24.9	41.7	10.7	7.5	9.5	6.9	91.8	92.4
<b>9330</b> survey area=3,414,940ha coverage=45,000ha (1.3%) known coverage=38,930ha (86.5%) ss mean=0.11 corr. presence/ss=0.57 M = 390, m = 46	1	46.0	14.9	127.6	46.0	15.7	47.7	10.0	87.6	67.3
	2	92.0	14.9	90.2	36.7	12.1	34.3	8.1	88.2	74.7
	3	138.0	14.7	73.7	32.9	10.8	28.2	7.1	87.2	76.8
	4	184.0	14.2	63.8	31.0	10.0	24.6	6.5	85.4	78.0
<b>9410</b> survey area=3,312,339ha coverage=260,000ha (7.9%) known coverage=210,884ha (81.1%) ss mean=0.14 corr. presence/ss=0.95 M = 368, m = 44	1	44.0	49.3	51.7	18.8	16.3	26.7	11.7	98.2	74.0
	2	88.0	49.4	36.5	13.9	11.6	19.1	9.6	98.6	81.7
	3	132.0	49.0	29.8	11.8	9.5	15.8	8.3	98.4	84.7
	4	176.0	48.1	25.8	10.6	8.3	13.7	7.5	98.3	86.7
<b>9420</b> survey area=3,874,989ha coverage=150,000ha (3.9%) known coverage=117,649ha (78.4%) ss mean=0.10 corr. presence/ss=0.91 M = 442, m = 51	1	51.0	39.7	69.8	21.3	17.1	25.7	13.3	96.4	77.1
	2	102.0	39.7	49.3	16.1	12.3	18.5	10.8	96.5	83.9
	3	153.0	39.3	40.3	13.8	10.1	15.2	9.2	96.1	86.8
	4	204.0	38.4	34.9	12.6	8.9	13.3	8.2	95.5	

ESS=expected sample size, EPS=expected presence in the sample (%), RSE=relative standard error (%), ERSEE=expectation of relative standard error estimator (%), C95=coverage of the 0.95 confidence interval (%), DE=double expansion estimator and DIF=difference estimator. Purple values refer to RSEs (%) achieved by the HT estimator under simple random sampling without replacement (SRSWOR) with ESS taken to have a fixed sample size. Values in blue and green refer to DE and DIF estimation, respectively.

the possibility of reaching them. For the remaining cells with  $h_j = 0$ , nonresponses were established by generating a random number  $u$  uniformly distributed in the interval  $(0, 1)$  and then assigning  $r_j = 0$  if  $u$  was smaller than a previously established nonresponse rate  $\rho$  and  $r_j = 1$  otherwise. To consider several levels of nonresponses, simulations were performed for  $\rho = 0.05, 0.10, 0.20$ . Once the  $r_j$ s were generated, they remained fixed throughout the simulation runs as fixed characteristics of the cells because of the design-based nature of the nonresponse treatment adopted in this study (see Section 7 of the Supporting Information

file).

For each habitat and each nonresponse level  $\rho$ , we independently performed  $R = 100,000$  two-phase selections of cells following the sampling scheme described in subSection 2.4. In the second phase, we selected a maximum number of  $\bar{n} = 4$  cells as the most sustainable effort. Because the number of nonrespondent cells in the final sample was a random variable, the expected number of nonresponses (ENR) was empirically computed as

$$ENR = \frac{1}{R} \sum_{r=1}^R nor_r$$

where  $nor_r$  is the number of nonresponses in the final sample selected at the  $r$ -th simulation run.

Moreover, for each selected sample, the estimators  $\hat{Y}_{(2)DE-CAL}$  and  $\hat{Y}_{(2)DIF-CAL}$  were computed from the sample data together with their variance estimators  $V_{DE-CAL}^2$  and  $V_{DIF-CAL}^2$ . For each habitat, each nonresponse rate and both estimators, we achieved the collection of the  $R$  coverage estimates,  $\hat{Y}_1, \dots, \hat{Y}_R$ , and the corresponding variance estimates,  $V_1^2, \dots, V_R^2$ , from which we derived the relative standard error estimates  $RSE_1, \dots, RSE_R$  with  $RSE_r = V_r/\hat{Y}_r$  for  $r = 1, \dots, R$ . Finally, the confidence interval at the nominal level of 0.95 was achieved by  $\hat{Y}_r \pm 2V_r$ . Then, from the resulting Monte Carlo distributions, we derived the expectation and the mean squared errors of the abundance estimator as follows:

$$E = \frac{1}{R} \sum_{r=1}^R \hat{Y}_r$$

and

$$MSE = \frac{1}{R} \sum_{r=1}^R (\hat{Y}_r - Y)^2.$$

From these quantities, the relative bias  $RB = (E - Y)/Y$  and the relative root mean squared error  $RRMSE = \sqrt{MSE}/Y$  were determined together with the expectation of the relative standard error estimators

$$ERSEE = \frac{1}{R} \sum_{r=1}^R RSE_r$$

and the actual coverage of the nominal 0.95 confidence intervals  $C95$  as the fraction of intervals containing the true coverage  $Y$ . Notably, neither estimator ensured design-unbiasedness. Therefore, their RB values were most important and their precision was quantified by their MSEs rather than by their variances.

### 3. Results

The values of relative bias (RB) concerning the double expansion (DE) and difference (DIF) estimators—theoretically equal to 0—were very close to 0 (always smaller than 0.7%).

The generated populations showed standardized HSS mean values of approximately 0.10 for most habitats except for habitat 6220 which showed approximately 0.03, and habitat 9410 which showed approximately 0.14. The way in which habitat presence was generated within cells gave rise to a strong correlation between  $x_j$ s and  $y_j$ s that varied between 0.80 and 0.97, except for habitat 9330, which showed a correlation of 0.57. Of the cells with habitat presence, the percentage of cells where habitat presence was known varied from 25% for habitat 6220 to 99% for habitat 9120 (Table 3).

The increase from 1 to 4 of the maximum number of cells to select in the second phase within the selected quadrats ( $\bar{n}$ ) led to considerable increases in the precision with RSEs of DE that halved their values in most cases. However, the reduction in RSEs decreases as  $\bar{n}$  increases, suggesting that further increments of  $\bar{n}$  over 4 are unsuitable, increasing the sampling effort without producing relevant improvements (Table 3 and Fig. 5).

The expected sample sizes (ESS) were invariably equal to  $\bar{n}m$ , showing that it was highly improbable to find quadrats where the number of c-blocks available for the second-phase sampling was smaller than  $\bar{n}$ . Of the total number of selected cells, the expected percentages of selected cells with habitat presence were much greater—from approximately 6 to approximately 40 times—than those expected under SRSWOR (Table 3).

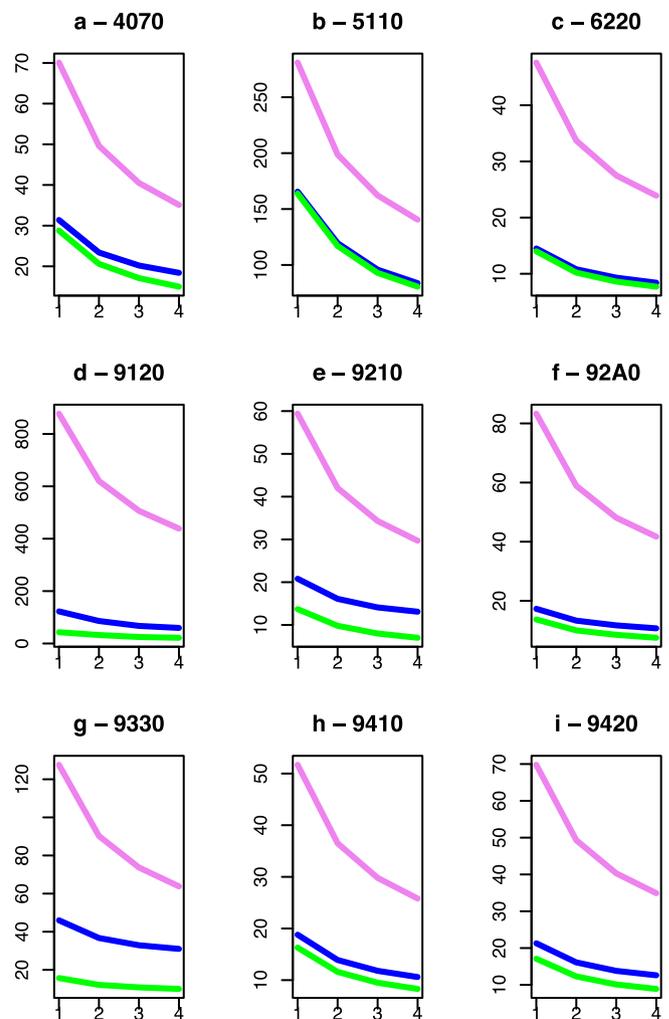


Fig. 5. Graphs of the relative standard error (RSE) in percentage plotted against  $\bar{n}$  from 1 to 4 for each habitat and for the Horvitz-Thompson estimator (HT) under simple random sampling without replacement (SRSWOR) and for the DE and DIF estimators under the proposed scheme (purple, blue and green lines, respectively, for the values in Table 3).

The precision of the DE estimator with respect to the crude HT estimator under SRSWOR was high, with efficiencies (ratio of RSEs) that varied from approximately 1.5 in habitat 5110 to approximately 7 in habitat 9120. In most cases, efficiencies ranged from 2 to 3 (Table 3 and Fig. 5).

The DIF estimator that exploited the previous knowledge of habitat presence in unsampled cells invariably outperformed the DE estimator (Table 3 and Fig. 5). Gains in precision were relevant when previous knowledge of habitat presence covered a large percentage, over 60%, of the habitat coverage (Fig. 5 a, d, e, g, h, i).

In absolute terms, for  $\bar{n} = 4$ , i.e., sampling fractions always smaller than 0.04%, the proposed sampling scheme combined with the DIF estimator yielded suitable precision with RSEs smaller than 15% when (i) populations were large (some millions of cells), (ii) HSSs were good proxies for habitat presence with correlations of  $x_j$ s vs  $y_j$ s of approximately 0.8–0.9, and (iii) the habitat coverage with respect to the survey areas was not smaller than 1%. These features were shared by habitats 4070, 6220, 9210, 92A0, 9410 and 9420. In these cases, RSEs were approximately 7–8%, except for habitat 4070, for which the RSE was 15% (Fig. 5 a, c, e, f, h, i). For habitats 5110, 9120 and 9330, the precision was unsuitable (Fig. 5 b, d, g). Particularly for habitat 5110, the RSE of the DE estimator was 83%, and the DIF estimator did not yield relevant improvements (Fig. 5 b).

The tentative RSE estimator that applied the HH criterion to the cells selected in the second phase using the product of the first-phase inclusion probabilities with those of the second phase as if they were the actual first-order inclusion probabilities provided unsatisfactory results, especially for the DIF estimator. While the proposed estimator performed quite well for the DE estimator, in most cases providing moderate overestimation with coverages of confidence intervals near the nominal level of 95%, it invariably underestimated the RSEs of the DIF estimator, with interval coverages invariably smaller than the nominal level. In some cases (e.g., habitat 9120), the underestimation was unsuitably large (Table 3). Therefore, as a precautionary rule of thumb, the RSE estimates achieved for the DE estimator should also be used to estimate the RSEs of the DIF estimator.

Since the RSEs of the DIF estimator were always smaller than those of the DE estimator, which in turn were overestimated, a fortiori these estimates should also overestimate the RSEs of the DIF estimator.

For the simulation study performed to check the nonresponse effects on the properties of the estimator, estimation based on the respondent sample was equivalent to estimation based on the complete sample when the sample values that could not be recorded were set to 0. Therefore, the effect of nonresponses turned out to be negligible if most of them occurred where the habitat was absent. Moreover, given that nonresponses were not allowed where the habitat presence was known from previous investigations, nonresponse effects were weak for habitats with high percentages of known coverages such as 9120, 9210 and 9330, where these percentages were greater than 85%. In these cases, even under a nonresponse rate of 20%, the negative bias was smaller than 3%, and the DIF estimator performed on the respondent samples provided the best results in terms of RRMSEs (Table 4).

In the other cases, when nonresponses also occurred where the

habitat was present, they heavily impacted the bias of both the DE and DIF estimators, with biases in some cases reaching levels of -15%. The bias that affected the estimators in the case of nonresponses decreased the precision, producing RRMSEs that were greater than the RSEs achieved with complete samples (Tables 3 and 4). At the same time, bias increased the underestimation of RRMSEs below the actual values and skewed the confidence intervals, decreasing their actual coverages. In these cases, calibration was necessary. However, the DE-CAL estimator eliminated negative bias at the cost of inducing a positive bias that sometimes was greater (in absolute value) than that entailed by nonresponses. On the other hand, the DIF-CAL estimator considerably reduced the negative bias even without reversing the sign and produced RRMSEs invariably smaller than those produced by DE-CAL. Therefore, the use of DIF-CAL seemed to be the best solution when calibration is necessary.

However, as in the case of complete samples, the HH-like variance estimator (SM.24) underestimated the actual precision, producing poor coverage of the resulting confidence intervals. In such cases, as a precautionary rule of thumb, the estimator (SM.22) achieved for the DE-CAL estimator should also be used to estimate the precision of the DIF-CAL estimator and to construct confidence intervals.

#### 4. Discussion

We developed an adaptive (sensu Lindenmayer and Likens, 2009, 2018; Lindenmayer et al., 2020) sampling strategy to provide unbiased estimators, or nearly unbiased in the presence of nonresponses, of habitat coverage over the study region (expressed as the number of 1ha cells occupied by habitat type) as an affordable and sound measure to satisfy the quantitative measurement of the “area” criterion in

**Table 4**

Monte-Carlo performance of double expansion and difference estimators of coverage calibrated to account for nonresponses and performance of relative standard error estimators for nine habitats in the study region and three nonresponse rates compared with the same indicators achieved from the respondent sample neglecting nonresponse presences

Habitat	NRR (%)	ENR	RESPONDENT SAMPLE								CALIBRATION							
			RB (%)		RRMSE (%)		ERSEE (%)		C95 (%)		RB (%)		RRMSE (%)		ERSEE (%)		C95 (%)	
			DE	DIF	DE	DIF	DE	DIF	DE	DIF	DE	DIF	DE	DIF	DE	DIF	DE	DIF
4070	5	5.1	-1.7	-1.7	18.3	19.9	20.6	17.9	94.0	91.2	2.4	-0.5	19.9	15.2	17.9	13.1	91.2	85.2
	10	10.2	-3.4	-3.4	18.4	14.8	20.7	15.5	92.9	81.9	4.4	-1.0	20.5	15.4	17.2	12.8	89.7	83.3
	20	20.2	-6.9	-6.9	19.0	15.2	20.9	13.1	90.3	75.4	8.8	-2.4	22.5	15.8	15.8	12.0	84.8	78.8
5110	5	2.8	-3.9	-3.9	81.8	78.5	62.8	59.0	81.3	68.1	1.4	-0.4	84.4	80.6	59.0	37.7	81.8	69.0
	10	5.5	-7.2	-7.2	80.5	77.2	63.4	38.5	79.8	65.8	3.1	-0.4	87.1	83.0	56.4	35.7	80.8	67.3
	20	11.0	-14.1	-14.2	78.0	74.3	64.4	37.0	76.6	60.9	6.6	-0.9	92.7	87.5	51.0	31.8	78.5	63.3
6220	5	28.6	-3.7	-3.7	9.1	8.4	9.0	8.1	92.4	90.1	1.2	-0.1	9.8	8.5	9.0	8.0	93.4	93.4
	10	57.1	-7.5	-7.5	11.1	10.5	9.3	8.2	84.5	80.9	2.2	-0.4	10.2	8.7	8.7	7.8	91.7	91.8
	20	114.2	-15.0	-15.05	17.0	16.6	9.7	8.4	55.8	46.0	4.7	-1.1	11.6	9.2	8.1	7.4	85.8	87.7
9120	5	0.5	-0.2	0.1	59.5	21.8	66.7	1.4	90.6	2.3	7.1	0.4	67.4	24.2	72.1	1.4	90.5	2.3
	10	1.0	-0.1	0.1	59.5	21.9	66.7	1.4	90.6	2.3	13.3	0.6	73.2	25.6	69.6	1.3	90.3	2.3
	20	2.0	-0.1	0.1	59.5	21.8	66.7	1.4	90.6	2.3	28.4	1.1	89.0	29.2	61.7	1.2	89.1	2.3
9210	5	8.05	-0.4	-0.4	13.0	6.8	14.6	4.9	96.2	67.3	3.1	-0.2	14.1	7.0	12.2	4.8	91.6	67.7
	10	16.1	-0.8	-0.8	13.0	6.7	14.6	4.7	96.1	65.1	6.2	-0.3	15.4	7.1	11.7	4.7	88.5	65.8
	20	32.2	-1.6	-1.6	12.9	6.4	14.6	4.4	95.7	60.3	13.1	-0.7	19.6	7.1	10.6	4.3	75.0	61.8
92A0	5	23.6	-2.0	-2.0	10.8	7.6	9.6	6.9	90.5	89.2	2.9	0.1	12.6	8.0	10.6	7.1	92.2	91.4
	10	46.9	-4.1	-4.1	11.2	8.2	9.7	6.7	87.9	84.0	5.9	-0.3	14.0	8.1	10.2	6.9	89.2	89.8
	20	93.4	-4.3	-4.3	12.9	9.1	13.4	7.6	92.3	76.0	12.5	-1.3	20.4	9.6	11.9	7.3	80.0	80.2
9330	5	8.1	-0.7	-0.6	30.9	9.7	24.7	6.4	85.1	75.9	4.6	0.1	33.7	10.3	24.5	6.4	85.0	76.9
	10	116.3	-1.4	-1.3	30.8	9.4	24.9	6.2	84.7	73.4	8.9	-0.0	36.0	10.4	23.4	6.2	83.7	75.0
	20	32.6	-2.8	-2.7	30.8	9.1	25.1	5.8	84.0	68.0	19.3	-0.2	42.7	10.7	21.2	5.7	77.9	71.3
4210	5	5.2	-0.9	-0.9	10.6	8.1	13.8	7.3	97.9	84.2	2.4	-0.4	12.3	8.3	11.6	7.1	93.1	85.0
	10	10.3	-1.9	-1.9	10.7	8.1	13.8	7.2	97.5	81.1	4.4	-0.8	13.0	8.4	11.1	6.9	90.5	82.7
	20	20.5	-3.8	-3.8	11.1	8.4	13.9	6.8	96.2	73.8	8.9	-1.8	15.4	8.6	10.2	6.5	81.4	77.7
9420	5	6.8	-1.0	-1.0	12.5	8.8	13.3	8.1	95.0	85.8	3.1	-0.2	15.3	9.2	13.6	8.0	82.4	86.5
	10	13.6	-2.2	-2.2	12.6	8.8	13.3	7.9	94.2	82.9	6.0	-0.6	16.4	9.3	13.0	7.8	80.0	80.2
	20	27.2	-4.3	-4.3	12.9	9.1	13.4	7.6	92.3	76.0	12.5	-1.3	20.4	9.6	11.9	7.3	80.0	80.2

NRR = nonresponse rate (%), ENR = expected nonresponses in the sample (%), RB = relative bias (%), RRMSE = relative root mean squared error (%), ERSEE = expectation of relative standard error estimator (%), C95 = coverage of the 95% confidence interval (%), DE = double expansion estimator, DIF = difference estimator. Values in blue and green refer to DE and DIF estimation, respectively.

compliance with the mandatory reporting cycle of the HD (ex art. 17). This approach can be adapted to any other country needing to soundly estimate the area covered by a habitat under study, vegetation type or ecosystem that cannot be properly mapped.

The results of a simulation study showed that the design-based inference performed by means of two-phase sampling and the use of the DIF estimator exploiting previous knowledge of habitat presence, or its calibrated counterpart in the presence of nonresponse have the potential to improve precision with respect to the HT estimator achieved under SRSWOR, thus showing a considerable design effect. Adopting a small sampling fraction not greater than 0.04% of the survey area, the DIF estimator provides suitable precision with RSEs smaller than 15% if habitats are quite common (e.g.,  $M > 200$ ), if the HSSs are good proxies for habitat presence and if coverages are not smaller than 1% with respect to the survey areas. Therefore, the same general strategy can be efficiently applied to different habitat types, from grassland (e.g., 6220) to forests (e.g., 9210), simply by changing the set of environmental predictors. Moreover, these habitats are characterized by a large geographic distribution across the three Italian biogeographical regions (Cervellini et al., 2021), partially confirming the applicability of the strategy to different macroecological and biogeographical contexts. On the other hand, the strategy provides unsuitable precision when the portion of coverage is small with respect to the survey area (e.g., 5110)—a characteristic that reduces the precision of any sampling strategy in spatial surveys—when the survey area is too small (e.g., 9120) and when the correlation between HSSs and habitat presence is weak (e.g., 9330). Notably, however, these situations are likely to reduce the precision of most sampling strategies and not only that of our strategy.

From these results, it is also apparent that for sufficiently large coverages, the correlation between HSSs and habitat presence and the information on habitat presence that may be available from previous investigations and surveys are determinant factors to efficiently increase the precision of the DIF estimator or its calibrated counterpart. A less satisfactory issue concerns the RSE estimation and the construction of confidence intervals. However, it is well known that this issue is “slightly tricky” in spatial sampling (e.g., Grafström, 2012), when, as in our case, the selection of neighbouring units is avoided or reduced.

These findings are essential for producing reliable estimates of area coverage by each of the 124 habitat types in the study region for the forthcoming fifth reporting cycle (2019–2024). In addition, for each habitat, our sampling strategy and the related estimation of its precision allow for a statistically sound detection of changes and trends (Lengyel et al., 2018) in relation to future national reports (e.g., 6th report). Furthermore, in the context of a “mandated” monitoring (Lindenmayer and Likens, 2010), it is now possible to plan a rigorous program based on a sustainable effort for each reporting cycle (Lindenmayer et al., 2020).

## 5. Concluding remarks

The design we developed is based on two-phase sampling that ensures the geographic spread of sample units in the area covered by each habitat type, permitting the collection of information across the whole habitat, as well as the use of local clusters of sample units that allow the optimization of time effort spent moving among them. This design permits the production of sound statistical estimates of habitat coverage while simultaneously providing sound information about uncertainty. In addition, exploitation of the difference estimator permits us to positively include all the habitat occurrence data that are collected out of the sampling scheme proposed here in the resulting estimates, such as data from local surveys or management or monitoring plans within protected areas. While the assemblage of these data, collected at local scales without a probabilistic design, does not facilitate statistically sound inferences, available data are effectively exploited in this approach as auxiliary components, improving the quality of the estimates produced on the basis of only the data collected by the sampling scheme. Basically,

this sampling strategy permits unification of the data collected under a well-designed probabilistic scheme with those opportunistically provided by all other available sources. Given the complex distribution of the various habitats, at the national, biogeographical and continental scales, this design-based approach permits the integration of a specifically defined and limited probabilistic sample with a likely larger sample lacking probabilistic features. Ultimately, this approach can be profitably used to arrange “mandated” monitoring plans at broad scales, such as the whole nation, biogeographical region or European Union, as is the case for the monitoring imposed by the Habitats Directive.

## Data availability statement

Data will be made available upon request.

## CRedit authorship contribution statement

**Lorenzo Fattorini:** Conceptualization. **Marco Cervellini:** Conceptualization. **Sara Franceschi:** Conceptualization, Data curation. **Michele Di Musciano:** Data curation. **Piero Zannini:** Data curation. **Alessandro Chiarucci:** Conceptualization.

## Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgements

Research contributing to this study was funded by the project “Development of a National Plan for Biodiversity Monitoring” (Italian National Institute for Environmental Protection and Research—ISPRA).

## Appendix A. Supplementary data

Supplementary data associated with this article can be found, in the online version, at <https://doi.org/10.1016/j.ecolind.2022.109352>.

## References

- Angelini, P., Casella, L., Grignetti, A. and Genovesi, P. (2016), Manuali per il monitoraggio di specie e habitat di interesse comunitario (Direttiva 92/43/CEE) in Italia: habitat, ISPRA, Serie Manuali e linee guida, 142/2016.
- Álvarez-Martínez, J.M., Jiménez-Alfaro, B., Barquín, J., Ondiviela, B., Recio, M., Silió-Calzada, A., Juanes, J.A., 2018. Modelling the area of occupancy of habitat types with remote sensing. *Methods Ecol. Evol.* 9 (3), 580–593. <https://doi.org/10.1111/2041-210X.12925>.
- Breidt, F.J., 1995. Markov chain designs for one-per-stratum sampling. *Survey Methodology* 21, 63–70.
- Brooks, T.M., Pimm, S.L., Akçakaya, H.R., Buchanan, G.M., Butchart, S.H.M., Foden, W., Hilton-Taylor, C., Hoffmann, M., Jenkins, C.N., Joppa, L., Rondinini, C., 2019. Measuring terrestrial area of habitat (AOH) and its utility for the IUCN Red List. *Trends Ecol. Evol.* 34 (11), 977–986. <https://doi.org/10.1016/j.tree.2019.06.009>.
- Brown, J.A., Robertson, B.L., McDonald, T., 2015. Spatially balanced sampling: application to environmental surveys. *Procedia Environ. Sci.* 27, 6–9. <https://doi.org/10.1016/j.proenv.2015.07.108>.
- Carli, E., Massimi, M., Angelini, P., Casella, L., Attorre, F., Agrillo, E., 2020. How to improve the distribution maps of habitat types at national scale. *Rendiconti Lincei. Scienze Fisiche e Naturali* 31 (3), 881–888. <https://doi.org/10.1007/s12210-020-00917-7>.
- CDR (2022). EIONET Central Data Repository. <https://cdr.eionet.europa.eu/it/eu/art17/envxuw6g/>. Accessed on June 21st 2022.
- Cervellini, M., Di Musciano, M., Zannini, P., Fattorini, S., Jiménez-Alfaro, B., Agrillo, E., Attorre, F., Angelini, P., Beierkuhnlein, C., Casella, L., Field, R., Fischer, J., Genovesi, P., Hoffmann, S., Irl, S.D.H., Nascimbene, J., Rocchini, D., Steinbauer, M., Vetaas, O.R., Chiarucci, A., 2021. Diversity of European habitat types is correlated with geography more than climate and human pressure. *Ecol. Evol.* <https://doi.org/10.1002/ece3.8409>.
- Cervellini, M., Zannini, P., Di Musciano, M., Fattorini, S., Jiménez-Alfaro, B., Rocchini, D., Field, R., Vetaas, O.R., Irl, S.D.H., Beierkuhnlein, C., Fischer, J.C., Casella, L., Angelini, P., Genovesi, P., Nascimbene, J., Chiarucci, A., 2020. A grid-

- based map for the Biogeographical Regions of Europe. Biodiversity Data J. 8 <https://doi.org/10.3897/BDJ.8.e53720>.
- Chiarucci, A., Di Biase, R.M., Fattorini, L., Marcheselli, M., Pisani, C., 2018. Joining the incompatible: Exploiting purposive lists for the sample-based estimation of species richness. *Ann. Appl. Stat.* 12 (3), 1679–1699. <https://doi.org/10.1214/17-AOAS1126>.
- CLC (2018), Corine Land Cover. Version is v.2020\_20u1. <https://land.copernicus.eu/pan-european/corine-land-cover/clc2018?tab=download>.
- Coops, N.C., Wulder, M.A., 2019. Breaking the Habit (at). *Trends Ecol. Evol.* 34 (7), 585–587. <https://doi.org/10.1016/j.tree.2019.04.013>.
- Cowles, J., Templeton, L., Battles, J.J., Edmunds, P.J., Carpenter, R.C., Carpenter, S.R., Paul Nelson, M., Cleavitt, N.L., Fahey, T.J., Groffman, P.M., Sullivan, J.H., Neel, M. C., Hansen, G.J.A., Hobbie, S., Holbrook, S.J., Kazanski, C.E., Seabloom, E.W., Schmitt, R.J., Stanley, E.H., Vander Zanden, J.M., 2021. Resilience: insights from the US LongTerm Ecological Research Network. *Ecosphere* 12 (5), e03434. <https://doi.org/10.1002/ecs2.3434>.
- Davies, C.E., Moss, D., Hill and M.O. (2004), EUNIS habitat classification revised 2004, Report to: European Environment Agency-European Topic Centre on Nature Protection and Biodiversity, 127–143.
- Delbos, P., Lagrange, I., Rozo, C., Bensettiti, F., Bouzillé, J.B., Evans, D., Lalanne, A., Rapinel, S., Bioret, F., 2021. Assessing the conservation status of coastal habitats under Article 17 of the EU Habitats Directive. *Biol. Conserv.* 254, 108935 <https://doi.org/10.1016/j.biocon.2020.108935>.
- DEM20, Digital Elevation Model. Rete del Sistema Informativo Nazionale Ambientale. SINAnet, <https://www.sinanet.isprambiente.it/it/sia-ispra/download-mais/dem20/view>.
- Drakou, E.G., Kallimanis, A.S., Mazaris, A.D., Apostolopoulou, E., Pantis, J.D., 2011. Habitat type richness associations with environmental variables: a case study in the Greek Natura 2000 aquatic ecosystems. *Biodivers. Conserv.* 20 (5), 929–943. <https://doi.org/10.1007/s10531-011-0005-4>.
- EIONET (2022), Article 17 web tool. Available from <https://nature-art17.eionet.europa.eu/article17/habitat/summary/?period=3&group=>.
- Ellwanger, G., Runge, S., Wagner, M., Ackermann, W., Neukirchen, M., Frederking, W., Müller, C., Ssymank, A., Sukopp, U., 2018. Current status of habitat monitoring in the European Union according to Article 17 of the Habitats Directive, with an emphasis on habitat structure and functions and on Germany. *Nature Conservation* 29 (October), 57–78. <https://doi.org/10.3897/natureconservation.29.27273>.
- European Commission (2020). Communication from the Commission to the European Parliament, the Council, the European Economic and Social Committee and the Committee of the Regions. EU Biodiversity Strategy for 2030 - Bringing nature back into our lives.
- European Commission - DG environment (2021). [https://ec.europa.eu/environment/natura2000/index\\_en.htm](https://ec.europa.eu/environment/natura2000/index_en.htm). Accessed July 2021 the 29th.
- European Environment Agency (2013), EEA Reference Grid. Available from <https://www.eea.europa.eu/data-and-maps/data/eea-reference-grids-2>.
- Fahrig, L., 2013. Rethinking patch size and isolation effects: the habitat amount hypothesis. *J. Biogeogr.* 40 (9), 1649–1663. <https://doi.org/10.1111/jbi.12130>.
- Fattorini, L., Franceschi, S., Maffei, D., 2013. Design-based treatment of unit nonresponse in environmental surveys using calibration weighting. *Biometrical J.* 55 (6), 925–943. <https://doi.org/10.1002/bimj.201100262>.
- Foster, S., 2020. DMBHdesign: An R-package for efficient spatial survey designs. *Methods Ecol. Evol.* 12 (3), 415–420. <https://doi.org/10.1111/2041-210X.13535>.
- Grafström, A., 2012. Spatially correlated Poisson sampling. *J. Stat. Planning Inference* 142 (1), 139–147. <https://doi.org/10.1016/j.jspi.2011.07.003>.
- Grafström, A., Lundström, N.L.P., 2013. Why well spread probability samples are balanced. *Open J. Stat.* 3 (1), 36–41. <https://doi.org/10.4236/ojs.2013.31005>.
- Hall, L.S., Krausman, P.R., Morrison, M.L., 1997. The Habitat Concept and a Plea for Standard Terminology. *Wildlife Society Bull.* 25 (1), 173–182.
- Hartigan, J.A., Wong, M.A., 1979. Algorithm AS 136: A k-means clustering algorithm. *J. R. Stat. Soc. Series C (Appl. Stat.)* 28 (1), 100–108. <https://doi.org/10.2307/2346830>.
- Haziza, D., Thompson, K.J., Yung, W., 2010. The effect of nonresponse adjustments on variance estimation. *Survey Methodol.* 36 (1), 35–43.
- ISTAT (2022), Confini delle unità amministrative a fini statistici, <https://www.istat.it/it/archivio/222527>. Accessed on June 21st 2022.
- ISTAT (2020), Sezioni di Censimento Litoranee. Shape file della linea litoranea, <https://www.istat.it/it/archivio/137341>. Accessed on March 2010.
- Keyes, A.A., McLaughlin, J.P., Barner, A.K., Dee, L.E., 2021. An ecological network approach to predict ecosystem service vulnerability to species losses. *Nature Commun.* 12 (1), 1–11. <https://doi.org/10.1038/s41467-021-21824-x>.
- Lengyel, S., Kosztyi, B., Schmeller, D.S., Henry, P.Y., Kotarac, M., Lin, Y.P., Henle, K., 2018. Evaluating and benchmarking biodiversity monitoring: Metadata-based indicators for sampling design, sampling effort and data analysis. *Ecol. Ind.* 85, 624–633. <https://doi.org/10.1016/j.ecolind.2017.11.012>.
- Lindenmayer, D.B., Likens, G.E., 2009. Adaptive monitoring: a new paradigm for long-term research and monitoring. *Trends Ecol. Evol.* 24 (9), 482–486. <https://doi.org/10.1016/j.tree.2009.03.005>.
- Lindenmayer, D.B., Likens, G.E., 2010. The science and application of ecological monitoring. *Biol. Conserv.* 143 (6), 1317–1328. <https://doi.org/10.1016/j.biocon.2010.02.013>.
- Lindenmayer, D.B., Likens, G.E., 2018. Maintaining the culture of ecology. *Front. Ecol. Environ.* 16 (4) <https://doi.org/10.1002/fee.1801>, 195–195.
- Lindenmayer, D.B., Woinarski, J., Legge, S., Southwell, D., Lavery, T., Robinson, N., Scheele, B., Wintle, B., 2020. A checklist of attributes for effective monitoring of threatened species and threatened ecosystems. *J. Environ. Manage.* 262, 110312 <https://doi.org/10.1016/j.jenvman.2020.110312>.
- Martinez-Harms, M.J., Bryan, A., Balvanera, P., Law, E.A., Rhodes, J.R., Possingham, H. P., Wilson, K.A., 2015. Making decisions for managing ecosystem services. *Biol. Conserv.* 184, 229–238. <https://doi.org/10.1016/j.biocon.2015.01.024>.
- Mitchell, S.C., 2005. How useful is the concept of habitat?—A critique. *Oikos* 110 (3), 634–638. <https://doi.org/10.1111/j.0030-1299.2005.13810.x>.
- MiTE (2021), National cartography of the network of sites. <https://www.mite.gov.it/pagina/cartografie-rete-natura-2000-e-aree-protette-progetto-natura>.
- Mulder, C., Bennett, E.M., Bohan, D.A., Bonkowski, M., Carpenter, S.R., Chalmers, R., Cramer, W., Durance, I., Eisenhauer, N., Fontaine, C., Houghton, A.J., Hettelingh, J. P., Hines, J., Ibanez, S., Jeppesen, E., Krumins, J.A., Ma, A., Mancinelli, G., Massol, F., Woodward, G., 2015. 10 years later: revisiting priorities for science and society a decade after the millennium ecosystem assessment. *Adv. Ecol. Res.* 53, 1–53. <https://doi.org/10.1016/bs.aecr.2015.10.005>.
- R Core Team (2020), R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria, <https://www.R-project.org/>.
- Sampford, M.R., 1967. On sampling without replacement with unequal probabilities of selection. *Biometrika* 54 (3–4), 499–513. <https://doi.org/10.2307/2335041>.
- Särndal, C.E., Swensson, B., Wretman, J., 1992. *Model assisted survey sampling*. Springer Science & Business Media.
- Smith, R.J., Gray, A.N., 2021. Strategic monitoring informs wilderness management and socioecological benefits. *Conservation Sci. Practice* 3 (9), e482. <https://doi.org/10.1111/csp.2482>.
- Wolter, K.M., 2007. *Introduction to variance estimation*, (2nd ed.). Springer, p. 53.
- SYapp, R.H., 1922. The concept of habitat. *J. Ecol.* 10 (1), 1–17.

## Supporting Information

for

# A sampling strategy for assessing habitat coverage at broad spatial scale

L. Fattorini, M. Cervellini, S. Franceschi, M. Di Musciano,

P. Zannini, A. Chiarucci

## 1 Statement of the problem

We suppose that the territory of interest, e.g., a national territory or an administrative district, is covered by a collection of quadrats of a given size, some of them lying partially outside the territory (see e.g., Figure 1). For any habitat we suppose to know the set  $Q$  of the  $M$  quadrats in which the habitat is present. Moreover, we partition each quadrat into  $K$  cells of an appropriate size suitable to be surveyed on the field for detecting habitat presence (e.g.,  $100m \times 100m$ ).

Henceforth, with a little abuse of notation, we identify quadrats and cells by their labels. Therefore, denote by  $V(q)$  the set of the  $K$  cells partitioning the quadrat  $q \in Q$ . Finally, we suppose that for each cell  $j \in V(q)$  and for each quadrat  $q \in Q$ , there is available a score  $x_j \geq 0$  indexing the likelihood of habitat presence in that cell. Scores can be determined by expert-based evaluation, models, remote sensing, or any other source of information. Henceforth, scores are referred to as habitat suitability score (HSS).

The HSSs will guide the sampling of cells: those cells for which  $x_j = 0$  have no possibility to host the habitat and as such they are excluded from the population to be sampled (e.g., the cells outside the territory of interest, those completely occupied by sea, etc.) while, as to the remaining cells, we will facilitate the selection of those cells with high HSS values. Therefore, denote by  $U(q) \subset V(q)$  the set of cells out of the  $K$  partitioning the quadrat  $q \in Q$  for which the HSSs take positive values, i.e., those for which there is a possibility to find the habitat. The set  $U = \cup_{q \in Q} U(q)$  is the population of cells to be sampled (see e.g., Figure 2) and  $N$  is the number of cells in  $U$ .  $N$  can be viewed as the extent (expressed in terms of the cell size) of the survey region, i.e., the region in which the habitat presence is possible.

For each cell  $j \in U$ , let  $y_j$  be the dichotomous quantity indexing habitat presence within the cell, i.e.,  $y_j = 1$  if the habitat is present and  $y_j = 0$ , otherwise. The parameter of interest to be estimated by means of a sample survey is the population total

$$Y = \sum_{j \in U} y_j$$

that represents the number of cells within the territory of interest in which habitat is present. The total  $Y$  can also be viewed as the habitat coverage, i.e., the extent (approximated at the cell size) in which the habitat is present.

For the subsequent developments, it is suitable to denote by

$$Y(q) = \sum_{j \in U(q)} y_j$$

the number of cells in the quadrat  $q \in Q$  in which the habitat is present, in such a way that  $Y$  can be rewritten as

$$Y = \sum_{q \in Q} Y(q)$$

Moreover, let

$$X(q) = \sum_{j \in U(q)} x_j$$

be the total of the HSSs within the quadrat  $q \in Q$ , that can be viewed as a suitability index for the whole quadrat.

We sampled the population  $U$  by means of a two-phase survey, with the first phase performed by two stages in which at the first stage a sample of quadrats is selected and then, in the second stage, a sample of cells is selected within each quadrat selected in the first stage, obtaining a first-phase sample of cells. Then, in the second phase, a sample of cells out of those selected in the first phase is selected in the final sample. At each stage of the first phase, the selection will be performed to ensure spatially balanced samples, i.e., samples of quadrats and cells well spread throughout the territory of interest as well as with probabilities increasing with the HSS values. In the second phase, only HSS values will guide the selection.

## 2 Sampling and estimation: first phase, first stage

In the first stage of the first phase, a sample of quadrats is selected adopting the one-per-stratum sampling (OPSS), to achieve a spatially balanced sample in which the selected quadrats are evenly spread throughout the territory, but also giving greater probability to the selection of those quadrats with greater HSS totals. To this purpose, the set  $Q$  of the  $M$  quadrats is partitioned into  $m$  blocks  $Q_1, \dots, Q_m$  of neighboring quadrats, henceforth referred to as q-blocks. The number of quadrats per q-block should be similar as possible, but not necessarily equal. That can be easily done, e.g., by using a clustering algorithm based on the spatial coordinates of the quadrat centroids.

Denote by

$$X(Q_l) = \sum_{q \in Q_l} X(q)$$

the HSS total in the  $l$ -th q-block. Then, a quadrat  $q \in Q_l$  is selected from the  $l$ -th q-block with probability proportional to its HSS total  $X(q)$ . Denote by  $q_l$  the label of the quadrat selected from the  $l$ -th q-block and by  $\Omega = \{q_1, \dots, q_m\}$  the sample of the  $m$  quadrats selected in the first stage.

The first-order inclusion probabilities of quadrats are given by

$$\omega_q = Pr\{q \in \Omega\} = \frac{X(q)}{X(Q_l)}, q \in Q_l, l = 1, \dots, m$$

while the-second order inclusion probabilities of two quadrats are given by

$$\omega_{qq'} = Pr \{q, q' \in \Omega\} = \begin{cases} \frac{X(q)X(q')}{X(Q_l)X(Q_{l'})}, & q \in Q_l, q' \in Q_{l'}, l \neq l' = 1, \dots, m \\ 0 & q, q' \in Q_l, l = 1, \dots, m \end{cases}$$

i.e., they are equal to the product of the first-order inclusion probabilities if the two quadrats belong to different q-blocks and are equal to 0 if the two quadrats belong to the same block.

If all the cells having positive HSS values in the  $m$  quadrats selected in the first stage of the first phase were visited and the habitat presences were recorded, then the totals  $Y(q_l)$  were known for each  $q_l \in \Omega$ . In this case, from the Horvitz-Thompson (HT) criterion, the first-phase, first stage unbiased estimator of  $Y$  would be

$$\hat{Y}_{(1,1)} = \sum_{l=1}^m \frac{X(Q_l)}{X(q_l)} Y(q_l) \quad (\text{SI.1})$$

From the independence of sampling performed within the  $m$  q-blocks, after little algebra, the variance of (SI.1) would be

$$Var_{\Omega} [\hat{Y}_{(1,1)}] = \sum_{l=1}^m X(Q_l) \sum_{q \in Q_l} \frac{Y^2(q)}{X(q)} - \sum_{l=1}^m Y^2(Q_l) \quad (\text{SI.2})$$

with

$$Y(Q_l) = \sum_{q \in Q_l} Y(q)$$

that is the number of cells in the  $l$ -th block  $Q_l$  in which the habitat is present, and the subscript  $\Omega$  denotes that uncertainty stems from all the possible samples of quadrats that can be selected by OPSS in the first stage of the first phase.

Regarding the number  $m$  of quadrats to be sampled, we use an increasing function  $m = m(M)$  such that the sampling fraction is equal to 1 when  $M$  is smaller than an established lower threshold  $M_L$ , decreases from 1 to an established minimum fraction  $f$  for  $M_L \leq M \leq M_U$ , and remains equal to  $f$  when  $M$  is greater than an established upper threshold  $M_U$ . Indeed, when  $M$  is small it is suitable and little expensive to visit all the  $M$  quadrats. In this case, the first stage is not performed, so that we pass directly to the second stage, i.e., the sampling of cells within the  $M$  quadrats. In other words, in this case there is no uncertainty induced by the first stage.

### 3 Sampling and estimation: first phase, second stage

Because in most cases it is prohibitive to visit all the cells with positive HSS values within the  $m$  quadrats selected in the first stage of the first phase, then the totals  $Y(q_l)$  within the selected quadrats are unknown. Therefore, the estimator (SI.1) is only virtual and it has been just reported to be used in the subsequent theoretical developments.

In practice, a second stage is necessary to estimate the totals  $Y(q_l)$  within the selected quadrats. To this purpose, a sample of cells is selected in each quadrat  $q_l \in \Omega$ , once again adopting OPSS to achieve spatially balanced samples of cells evenly spread throughout the selected quadrats. To

perform OPSS, each quadrat  $q_l \in \Omega$  constituted by the set  $V(q_l)$  of  $K$  cells is partitioned into  $k$  square blocks  $V(q_l, 1), \dots, V(q_l, k)$  of  $K/k$  adjacent cells (see e.g., Figure 3), henceforth referred to as c-blocks, with  $K$  multiple of  $k$ . Denote by  $X(q_l, i)$  the total of the HSSs within the  $i$ -th c-block ( $i = 1, \dots, k$ ) in the quadrat  $q_l \in \Omega$ . The c-blocks for which  $X(q_l, i) = 0$ , i.e., those that do not contain cells with positive HSSs, are discarded from the second stage of sampling. Denote by  $C(q_l)$  the set of labels identifying the  $k(q_l)$  c-blocks, out of the  $k$  partitioning the quadrat  $q_l \in \Omega$ , such that  $X(q_l, i) > 0$ , i.e., those c-blocks containing at least one cell with a positive HSS. Then, one cell is selected in each of these  $k(q_l)$  c-blocks with probability proportional to HSSs, excluding those cells with HSSs equal to 0. Denote by  $U(q_l, i)$  the set of cells in the c-block  $V(q_l, i)$  with  $i \in C(q_l)$  that have positive HSSs and denote by  $j(q_l, i)$  the label of the cell selected from  $U(q_l, i)$ . Therefore,  $\Theta(q_l) = \{j(q_l, i) : i \in C(q_l)\}$  is the sample of the  $k(q_l) \leq k$  cells selected in the quadrat  $q_l \in \Omega$ . Then,  $\Theta = \bigcup_{l=1}^m \Theta(q_l)$  is the final sample of cells selected at the end of the second stage of the first phase. Obviously, the maximum size of  $\Theta$  is  $m \times k$ , that is achieved when all the  $k$  c-blocks partitioning the  $m$  quadrats selected in the first stage contain cells with positive HSSs.

The first-order inclusion probabilities of the cells within the selected quadrat  $q_l \in \Omega$  are

$$\theta_j = Pr \{j \in \Theta(q_l) \vee \Omega\} = \frac{x_j}{X(q_l, i)}, j \in U(q_l, i), i \in C(q_l)$$

while the second-order inclusion probabilities of two cells within the same quadrat  $q_l \in \Omega$  are

$$\theta_{jj'} = Pr \{j, j' \in \Theta(q_l) \vee \Omega\} = \begin{cases} \frac{x_j}{X(q_l, i)} \frac{x_{j'}}{X(q_l, i')} & j \in U(q_l, i), j' \in U(q_l, i'), i \neq i' \in C(q_l) \\ 0 & j, j' \in U(q_l, i), i \in C(q_l) \end{cases}$$

i.e., they are equal to the product of the first-order inclusion probabilities if the two cells belong to different c-blocks of  $q_l$  and are equal to 0 if the two cells belong to the same c-block of  $q_l$ .

If all the cells selected by OPSS in the second stage within each quadrat  $q_l \in \Omega$  were visited and the habitat presence was recorded, then the totals  $Y(q_l)$  would be unbiasedly estimated for each  $q_l \in \Omega$  adopting the HT criterion by means of

$$\hat{Y}_{(1,2)}(q_l) = \sum_{i \in C(q_l)} X(q_l, i) \frac{y_{j(q_l, i)}}{x_{j(q_l, i)}} \quad (\text{SI.3})$$

Therefore, substituting (SI.3) into (SI.1), the first-phase unbiased estimator of  $Y$  would be

$$\hat{Y}_{(1)} = \sum_{l=1}^m \frac{X(Q_l)}{X(q_l)} \hat{Y}_{(1,2)}(q_l) \quad (\text{SI.4})$$

From the independence of estimation performed within c-blocks, after little algebra, the variance of (SI.3) would be

$$Var_{\Theta} \left[ \hat{Y}_{(1,2)}(q_l) \mid \Omega \right] = \sum_{i \in C(q_l)} X(q_l, i) \sum_{j \in U(q_l, i)} \frac{y_j}{x_j} - \sum_{i \in C(q_l)} Y^2(q_l, i) \quad (\text{SI.5})$$

where  $Y(q_l, i)$  is the number of cells within the c-block  $i \in C(q_l)$  containing the habitat and the sub-fix  $\Theta$  denotes that uncertainty stems from all the samples of cells that can be selected in the second stage by OPSS, conditional to the sample of quadrats  $\Omega$  selected in the first stage.

Regarding the variance of the first-phase estimator (SI.4), from the unbiasedness of the estimators (SI.3) and from the expression of their conditional variance (SI.5), it follows that

$$\begin{aligned} Var_{\Omega \circ \Theta} [\hat{Y}_{(1)}] &= Var_{\Omega} \left[ \sum_{l=1}^m \frac{X(Q_l)}{X(q_l)} E_{\Theta} \left\{ \hat{Y}_{(1,2)}(q_l) | \Omega \right\} \right] + E_{\Omega} \left[ \sum_{l=1}^m \frac{X^2(Q_l)}{X^2(q_l)} Var_{\Theta} \left\{ \hat{Y}_{(1,2)}(q_l) | \Omega \right\} \right] \\ &= Var_{\Omega} [\hat{Y}_{(1,1)}] + E_{\Omega} \left[ \sum_{l=1}^m \frac{X^2(Q_l)}{X^2(q_l)} Var_{\Theta} \left\{ \hat{Y}_{(1,2)}(q_l) | \Omega \right\} \right] \end{aligned} \quad (SI.6)$$

where the first term is given by equation (SI.2) and the sub-fix  $\Omega \circ \Theta$  denotes that uncertainty stems from all the samples of cells that can be selected at the end of the two stages that compose the first phase. Obviously, if the first stage is not performed, the first term of equation (SI.6) disappears and the quantity within square brackets in the second term is a constant rather than a random variable that depends on  $\Omega$ .

## 4 Sampling and estimation: second phase

In most cases it is also prohibitive to visit all the cells selected in the second stage within the  $m$  quadrats selected in the first stage of the first phase. Therefore, the estimators of totals within selected quadrats are only virtual, in such a way that also the first-phase estimator (SI.5) is virtual and its equation is just reported to be used in the subsequent theoretical development.

To reduce sampling efforts, a second phase is performed. If  $k(q_l)$  is equal or smaller than a given threshold  $\bar{n} < k$ , then all the cells in the sample  $\Theta(q_l)$  are sampled and visited. In this case, the second phase of sampling does not take place in the quadrat  $q_l$  and its total  $Y(q_l)$  is estimated by means of (SI.3). On the other hand, if  $k(q_l) > \bar{n}$ , then a sample  $\Pi(q_l)$  of  $\bar{n}$  cells is selected from  $\Theta(q_l)$  in the second phase by means of a fixed-size without-replacement sampling scheme with first-order inclusion probabilities increasing with HSSs. Therefore  $\Pi = \bigcup_{l=1}^m \Pi(q_l)$  is the final sample of cells selected at the end of the second phase, with  $\Pi(q_l) = \Theta(q_l)$  when  $k(q_l) \leq \bar{n}$ . Denote by  $n$  the random size of the final sample  $\Pi$ . Obviously, the maximum size of  $\Pi$  is  $m \times \bar{n}$ , that occurs when in all the quadrats  $q_l \in \Omega$  selected in the first stage of the first phase, the number of cells selected in the second stage of the first phase is greater than  $\bar{n}$ .

There is a long list of fixed-size schemes available to select samples with first-order inclusion probabilities increasing with an auxiliary variable. Hedayat and Sinha (1991, Chapter 5) devote an entire chapter to the use of auxiliary size measures (HSS in our case) in sampling, also arguing about the theoretical appealing of these schemes when the size measure is correlated with the survey variable. These schemes are partitioned into two broad classes: i) schemes with selection probabilities proportional to HSSs, usually referred to as PPS schemes; ii) schemes with inclusion probabilities proportional to HSSs, usually referred to as IIPS schemes. Both classes have strengths and weaknesses. We have opted for the more familiar IIPS schemes.

Regarding the IIPS schemes, the first-order inclusion probabilities proportional to the  $x_j$ s are given by

$$\pi_j = Pr \{j \in \Pi(q_l) | \Omega, \Theta\} = \frac{\bar{n}x_j}{X\{\Theta(q_l)\}}, j \in \Theta(q_l) \quad (SI.7)$$

where  $X\{\Theta(q_l)\}$  is the total of the  $x_j$ s in the sample of cells selected within  $q_l$  at the end of the first phase, ensuring that the sum of  $\pi_j$ s equals the sample size  $\bar{n}$ . Sometimes it could be possible that  $x_j > X\{\Theta(q_l)\}/\bar{n}$  for some  $j \in \Theta(q_l)$ , in such a way that  $\pi_j$  would be greater than 1. In these cases, the  $\pi_j$ s are set to be 1, i.e., these cells are selected with certainty and the PIPS selection is performed, *mutatis mutandis*, for the remaining cells of  $\Theta(q_l)$ .

While in PIPS schemes the first-order inclusion probabilities are very simple to compute, the sampling schemes ensuring that cells are selected in accordance with these probabilities are instead quite complex. Brewer and Hanif (1983) list 43 sampling schemes giving rise to first-order inclusion probabilities of type (SI.7), while more recent proposals are given by Chen et al. (1994), Tillé (1996) and Deville and Tillé (1998). All these schemes differ in the second-order inclusion probabilities, that have huge expressions in most cases. Among these schemes, we have opted for the Sampford scheme (Sampford, 1967) that is one of the most familiar PIPS scheme and has been automatized in many sampling software. The Sampford scheme is rejective, in the sense that it performs  $\bar{n}$  with replacement selections and if the resulting sample is not constituted by distinct units it is rejected until a sample of  $\bar{n}$  distinct units is achieved.

For simplicity of notation, it should be pointed out that in the second stage of the first phase a unique cell is selected within each of the  $k(q_l)$  c-blocks of  $q_l$ , in such a way that the  $\bar{n}$  cells selected in the second phase, i.e., those in the sample  $\Pi(q_l) \subset \Theta(q_l)$ , can be univocally determined by the labels of the c-blocks they belong. Therefore, by a little abuse of notation, we can simply denote by  $i \in \Pi(q_l)$  the cell  $j(q_l, i)$  selected in the second phase out of the  $k(q_l)$  cells of  $\Theta(q_l)$  selected at the end of the first phase. Accordingly, once the second phase is performed, the second-phase estimator of  $Y(q_l)$  can be written as

$$\hat{Y}_{(2)}(q_l) = \sum_{i \in \Pi(q_l)} X(q_l, i) \frac{y_{j(q_l, i)}}{x_{j(q_l, i)} \pi_{j(q_l, i)}} \quad (\text{SI.8})$$

Therefore, adopting (SI.8) into the first-phase estimator (SI.4), the second-phase estimator of  $Y$  is finally given by

$$\hat{Y}_{(2)DE} = \sum_{l=1}^m \frac{X(Q_l)}{X(q_l)} \hat{Y}_{(2)}(q_l) \quad (\text{SI.9})$$

Estimator (SI.9) is unbiased and is usually referred to as the double expansion (DE) estimator, being the results of the double application of the HT criterion (e.g., Särndal et al. 1992, section 9.3). From the well-known results of HT estimation (e.g., Hedayat and Sinha, 1991, Section 3.1), the variance of (SI.8) is given by

$$\begin{aligned} & \text{Var}_{\Pi} \left[ \hat{Y}_{(2)}(q_l) \vee \Omega, \Theta \right] \\ &= \sum_{h > i \in C(q_l)} (\pi_{j(q_l, i)} | \pi_{j(q_l, h)} - \pi_{j(q_l, i), j(q_l, h)}) \left[ \frac{X(q_l, i) y_{j(q_l, i)}}{x_{j(q_l, i)} \pi_{j(q_l, i)}} - \frac{X(q_l, h) y_{j(q_l, h)}}{x_{j(q_l, h)} \pi_{j(q_l, h)}} \right]^2 \end{aligned} \quad (\text{SI.10})$$

where  $\pi_{j(q_l, i), j(q_l, h)}$  denotes the second-order inclusion probability of the cells  $j(q_l, i)$  and  $j(q_l, h)$  and the sub-fix  $\Pi$  denotes that uncertainty stems from all the samples of cells that can be selected in the second phase by the Sampford scheme, conditional to the quadrats and cells selected in the first phase.

Regarding the variance of the second-phase estimator (SI.9), from the unbiasedness of  $\hat{Y}_{(2)}(q_l)$  and from the expression of the variances (SI.2), (SI.6) and (SI.10), it follows that

$$\begin{aligned}
& \text{Var}_{\Omega \circ \Theta \circ \Pi} \left[ \hat{Y}_{(2)DE} \right] \\
&= \text{Var}_{\Omega \circ \Theta} \left[ \sum_{l=1}^m \frac{X(Q_l)}{X(q_l)} E_{\Pi} \left\{ \hat{Y}_{(2)}(q_l) \mid \Omega, \Theta \right\} \right] + E_{\Omega \circ \Theta} \left[ \sum_{l=1}^m \frac{X^2(Q_l)}{X^2(q_l)} \text{Var}_{\Pi} \left\{ \hat{Y}_{(2)}(q_l) \mid \Omega, \Theta \right\} \right] \\
&= \text{Var}_{\Omega \circ \Theta} \left[ \sum_{l=1}^m \frac{X(Q_l)}{X(q_l)} \hat{Y}_{(1,2)}(q_l) \right] + E_{\Omega \circ \Theta} \left[ \sum_{l=1}^m \frac{X^2(Q_l)}{X^2(q_l)} \text{Var}_{\Pi} \left\{ \hat{Y}_{(2)}(q_l) \mid \Omega, \Theta \right\} \right] \\
&= \text{Var}_{\Omega \circ \Theta} \left[ \hat{Y}_{(1)} \right] + E_{\Omega \circ \Theta} \left[ \sum_{l=1}^m \frac{X^2(Q_l)}{X^2(q_l)} \text{Var}_{\Pi} \left\{ \hat{Y}_{(2)}(q_l) \mid \Omega, \Theta \right\} \right] \\
&= \text{Var}_{\Omega} \left[ \hat{Y}_{(1,1)} \right] + E_{\Omega} \left[ \sum_{l=1}^m \frac{X^2(Q_l)}{X^2(q_l)} \text{Var}_{\Theta} \left\{ \hat{Y}_{(1,2)}(q_l) \mid \Omega \right\} \right] + E_{\Omega \circ \Theta} \left[ \sum_{l=1}^m \frac{X^2(Q_l)}{X^2(q_l)} \text{Var}_{\Pi} \left\{ \hat{Y}_{(2)}(q_l) \mid \Omega, \Theta \right\} \right] \\
&= \sum_{l=1}^m \left[ X(Q_l) \sum_{q \in Q_l} \frac{Y^2(q)}{X(q)} - Y^2(Q_l) \right] + E_{\Omega} \left[ \sum_{l=1}^m \frac{X^2(Q_l)}{X^2(q_l)} \text{Var}_{\Theta} \left\{ \hat{Y}_{(1,2)}(q_l) \mid \Omega \right\} \right] \\
&\quad + E_{\Omega \circ \Theta} \left[ \sum_{l=1}^m \frac{X^2(Q_l)}{X^2(q_l)} \text{Var}_{\Pi} \left\{ \hat{Y}_{(2)}(q_l) \mid \Omega, \Theta \right\} \right] \tag{SI.11}
\end{aligned}$$

where the sub-fix  $\Omega \circ \Theta \circ \Pi$  denotes that uncertainty stems from all the samples of cells that can be selected at the end of the two phases. Obviously, when the first stage in the first phase is not performed, i.e., when  $M = m$ , or the second phase within a c-block is not performed, i.e., when  $k(q_l) \leq \bar{n}$ , the corresponding variance components vanish.

## 5 Precision estimation

To estimate the overall variance (SI.11), it should be considered that the OPSS schemes adopted in the first phase to select quadrats and cells involve second-order inclusion probabilities not invariably positive. Therefore, there is no possibility to estimate unbiasedly the first two terms of (SI.11), even if the third term could be estimated unbiasedly by means of the Sen-Yates-Grundy estimator (e.g., Hedayat and Sinha, 1991, section 3.1). Accordingly, *ad hoc* variance estimators are necessary. As is customary in similar situations, we have exploited the Hansen-Hurvitz (HH) criterion (e.g., Wolter, 2007), by presuming independence between the  $n$  units selected in the final sample. Because the first-order inclusion probabilities of the two-phase sampling scheme here adopted are unknown, we modify the HH variance estimator adopting the product of the first-order inclusion probabilities of

the first phase with those of the second phase. The resulting HH-like variance estimator turns out to be

$$V_{DE}^2 = \frac{1}{n(n-1)} \sum_{l=1}^m \sum_{i \in \Pi(q_l)} \left[ t_{j(q_l, i)} - \hat{Y}_{(2)DE} \right]^2 \quad (\text{SI.12})$$

with

$$t_{j(q_l, i)} = n \frac{y_{j(q_l, i)}}{\tau_{j(q_l, i)}}$$

and

$$\tau_{j(q_l, i)} = \frac{X(q_l)}{X(Q_l)} \frac{x_{j(q_l, i)}}{X(q_l, i)} \pi_{j(q_l, i)}$$

is the product of the first-phase with the second-phase first-order inclusion probabilities.

## 6 Exploiting auxiliary information by difference estimation

From a technical point of view, in survey sampling auxiliary information is constituted by variables correlated with the survey variable whose values are freely known or known at low costs for all the units in the population. Owing to the relationship between auxiliary and survey variable, the auxiliary information can be exploited by several criteria to improve estimation without increasing the sampling effort.

In this framework, owing to previous investigations performed by naturalists, there may exist a subset  $U_h \subset U$  of cells in which the habitat presence is claimed. Therefore, in accordance with the difference (DIF) estimation criterion (Särndal et al. 1992, Section 6.3), for each cell  $j$  in the population we can exploit a proxy for  $y_j$ , say  $h_j$ , that is known for each population unit and is equal to 1 for each  $j \in U_h$  and is equal to 0 otherwise. Obviously, the total of the  $h_j$ s

$$H = \sum_{j \in H} h_j$$

constitutes the extent in which the habitat presence is sure. Accordingly, the habitat coverage can be rewritten as

$$Y = H + \sum_{j \in U} d_j = H + D$$

where  $d_j = y_j - h_j$  is the error achieved by predicting  $y_j$  by means of its proxy  $h_j$ . In practice, the  $d_j$ s are equal to 0 for all the cells in which the habitat is absent as well as for all the cells in which the presence is known, while it is equal to 1 for all the cells in which the presence occurs but is not known.

While  $H$  is known, the total of errors  $D$  is obviously unknown and as such it can be estimated by the same strategy previously adopted to estimate  $Y$ . In other words, under the strategy delineated in sections 1-4, the difference estimator of  $Y$  becomes

$$\hat{Y}_{(2)DIF} = H + \hat{D}_{(2)DE} = H + \sum_{l=1}^m \frac{X(Q_l)}{X(q_l)} \hat{D}_{(2)}(q_l) \quad (\text{SI.13})$$

with

$$\hat{D}_{(2)}(q_l) = \sum_{i \in C(q_l)} X(q_l, i) \frac{d_{j(q_l, i)}}{x_{j(q_l, i)} \pi_{j(q_l, i)}}$$

In practice, in analogy to (SI.9),  $\hat{D}_{(2)DE}$  is the DE estimator of  $D$ . Because the  $d_j$ s are invariably nonnegative,  $\hat{D}_{(2)DE}$  is also nonnegative, in such a way that  $\hat{Y}_{(2)DIF}$  never provides inconsistent results, as it invariably exceeds the number of cells of sure presence  $H$ .

From the properties of the two-phase estimator,  $\hat{D}_{(2)DE}$  is an unbiased estimator of  $D$ , in such a way that  $\hat{Y}_{(2)DIF}$  is an unbiased estimator of  $Y$ . Moreover, owing to the presence of the constant term  $H$ , the variance of  $\hat{Y}_{(2)DIF}$  coincides with the variance of  $\hat{D}_{(2)DE}$  that, in analogy with equation (SI.11), is given by

$$\begin{aligned} \text{Var}_{\Omega \circ \Theta \circ \Pi} [\hat{D}_{(2)DE}] &= \sum_{l=1}^m \left[ X(Q_l) \sum_{q \in Q_l} \frac{D^2(q)}{X(q)} - D^2(Q_l) \right] \\ &+ E_{\Omega} \left[ \sum_{l=1}^m \frac{X^2(Q_l)}{X^2(q_l)} \text{Var}_{\Theta} \left\{ \hat{D}_{(1,2)}(q_l) \mid \Omega \right\} \right] + E_{\Omega \circ \Theta} \left[ \sum_{l=1}^m \frac{X^2(Q_l)}{X^2(q_l)} \text{Var}_{\Pi} \left\{ \hat{D}_{(2)}(q_l) \mid \Omega, \Theta \right\} \right] \end{aligned} \quad (\text{SI.14})$$

where  $D(Q_l)$ ,  $D(q_l)$  and  $D(q_l, i)$  are the totals of the  $d_j$ s within the  $q$ -block  $Q_l$ , the quadrat  $q_l$  and the  $i$ -th  $c$ -block of  $q_l$ , respectively,  $\hat{D}_{(1,2)}(q_l)$  is the second-stage estimator of  $D(q_l)$  and

$$\text{Var}_{\Theta} \left\{ \hat{D}_{(1,2)}(q_l) \mid \Omega \right\} = \sum_{i \in C(q_l)} X(q_l, i) \sum_{j \in U(q_l, i)} \frac{d_j}{x_j} - \sum_{i \in C(q_l)} D^2(q_l, i)$$

while

$$\begin{aligned} &\text{Var}_{\Pi} \left[ \hat{D}_{(2)}(q_l) \mid \Omega, \Theta \right] \\ &= \sum_{h > i \in C(q_l)} \left( \pi_{j(q_l, i)} \mid \pi_{j(q_l, h)} - \pi_{j(q_l, i), j(q_l, h)} \right) \left[ \frac{X(q_l, i) d_{j(q_l, i)}}{x_{j(q_l, i)} \pi_{j(q_l, i)}} - \frac{X(q_l, h) d_{j(q_l, h)}}{x_{j(q_l, h)} \pi_{j(q_l, h)}} \right]^2 \end{aligned}$$

Therefore, in analogy with equation (SI.12), the HH-like variance estimator is given by

$$V_{DIF}^2 = \frac{1}{n(n-1)} \sum_{l=1}^m \sum_{i \in \Pi(q_l)} \left[ t_{j(q_l, i)} - \hat{D}_{(2)DE} \right]^2 \quad (\text{SI.15})$$

where, in this case,

$$t_{j(q_l, i)} = n \frac{d_{j(q_l, i)}}{\tau_{j(q_l, i)}}$$

Even if the HSSs have been massively exploited at design level by selecting quadrats and cells with probabilities increasing with the HSS values, in principle we could further exploit the information provided by HSS also at the estimation level by calibrating the DE estimator by means of the ratio (RAT) estimation criterion (e.g., Hedayat and Sinha, 1991, Section 6.1, Särndal et al. 1992, Section 5.6). However, owing to the exploitation of HSS at design level, the subsequent exploitation of the same information at estimation level is likely to provide poor improvement with respect to the DE and DIF estimators. For this reason, that has been confirmed also from previous simulation studies, we do not pursue the use of RAT estimator. Rather, because HSSs should constitute a suitable auxiliary information available for all the cells in the population, we have used the RAT estimator to correct the bias induced by the presence of nonresponse.

## 7 Nonresponse treatment by calibration weighting

Unit nonresponse is often a problem in environmental surveys. It usually occurs when a population of spatial units partitioning a study area (the cells in our case) is sampled and some sampled units are in difficult terrains or are bounded by artificial barriers and cannot be reached by the survey crew, or, even if the unit is reached, the steep slope of the terrain does not allow the recording of the survey variable within. If the survey variable is nonnegative as in our case, any estimator performed on the respondent sample will be negatively biased. Therefore, a vast literature deals with the problem of nonresponse adjustment by means of suitable estimation techniques for reducing nonresponse bias.

Widely applied methods were referred to as nonresponse propensity weighting by Haziza et al. (2010). These methods view the respondent set as the result of a last phase of sampling. In the previous phases, a sample is selected from the population by means of the established sampling scheme. Then, in the last phase, it is assumed the existence of a response mechanism for which every population unit has its own invariably positive response probability. It is also assumed that each unit responds independently to the others. Practically speaking, the respondent set is viewed as a subsample of the selected sample, achieved using a sort of Poisson sampling. A realistic model is adopted to link the unknown response probabilities with some auxiliary variables. The model is subsequently used to estimate the response probabilities based on auxiliary information available, while the values of the survey variable are treated as fixed quantities, as customary in design-based inference. However, as pointed out by Fattorini et al. (2013), in environmental surveys unit responses cannot be viewed as outcomes of dichotomous and independent experiments with unknown probabilities. For example, in our framework, if some selected cells cannot be reached, no random experiment can be claimed because they will never be reached. In these situations, responses should be viewed as fixed characteristics of the cells, in such a way that the population is partitioned into respondent and nonrespondent strata, i.e., cells that can or cannot be reached. Moreover, the assumption that cell responses are independent events is likely to be unrealistic in our framework. Neighbouring cells, lying in terrains with the same characteristics, tend to have a similar response pattern, i.e., a sort of spatial contagion is likely to be present among responses. Based on these considerations, the use of nonresponse propensity weighting in environmental surveys, and in our surveys as well, does not seem to be logically defensible.

Alternatively, unit nonresponse could be handled by a plethora of imputation techniques, i.e., procedures in which nonresponse values are replaced by substitutes and estimation is performed

on the completed data, achieving the so-called imputed estimator. As pointed out by Särndal and Lundström (2005, p. 52), imputed values are artificial and are customarily obtained by means of a prediction model presuming a relationship between the survey variable and a set of auxiliary variables. In accordance with the presumed model, commonly used techniques of imputation are, among others, regression imputation, nearest-neighbour imputation, hot deck imputation, and multiple imputation. For a review see, for example, Little and Rubin (2002). Without entering on these techniques, it should be pointed out that, irrespective of the model adopted to make imputation, it cannot be validated in the set of nonrespondent units. Thus, in our opinion, it is difficult to scientifically defend any proposed imputation estimator.

Because response and prediction modelling are not convincing for treating nonresponse in our framework, we here follow the proposal by Fattorini et al. (2013) that adopt a complete design-based treatment in which population values and nonresponse are both viewed as fixed characteristics. To this purpose, the authors exploit the nonresponse calibration weighting by Haziza et al. (2010), that modifies the weights originally attached to each sampled units in such a way that new weights can estimate the population totals of a set of auxiliary variables without error. The rationale behind nonresponse calibration weighting is obvious: if the calibrated weights guess the totals of the auxiliary variables without errors, they should be suitable also for estimating the totals of the survey variable, providing that a relationship exists between survey and auxiliary variables. Moreover, if a perfect linear relationship exists between the survey variable and the auxiliary variables in the whole population, the nonresponse calibration weighting estimates the population total without error (Särndal and Lundström, 2005, p. 61). Accordingly, even if a perfect linear relationship does never hold in practice, nonresponse calibration weighting is likely to perform well for reducing nonresponse bias under strong linear relationships between survey and auxiliary variables.

Turning to our problem of estimating the habitat coverage  $Y$ , denote by  $U_R$  the set of cells in the population that can be reached by naturalist crews and  $U - U_R$  the set of cells that cannot be reached. Moreover, denote by  $z_j$  a response indicator variable that is equal to 1 if  $j \in U_R$  and equal to 0, otherwise, and denote by  $\Pi_R \subset \Pi$  the respondent sample, i.e., the set of cells in the final sample  $\Pi$  that have been reached and surveyed to record the habitat presence. The DE estimator (SI.9) based on the complete sample  $\Pi$  can be rewritten as a weighted sum of the sample observations, i.e.

$$\hat{Y}_{(2)DE} = \sum_{j \in \Pi} w_j y_j \quad (\text{SI.16})$$

it is at once apparent that the same estimator performed on the response sample  $\Pi_R$ , say

$$\hat{Y}_{(2)R} = \sum_{j \in \Pi_R} w_j y_j \quad (\text{SI.17})$$

turns out to be invariably smaller than  $\hat{Y}_{(2)DE}$ , i.e.

$$\hat{Y}_{(2)R} = \sum_{j \in \Pi_R} w_j y_j \leq \sum_{j \in \Pi} w_j y_j = \hat{Y}_{(2)DE}$$

from which  $E_{\Omega \circ \Theta \circ \Pi} \left\{ \hat{Y}_{(2)R} \right\} \leq E_{\Omega \circ \Theta \circ \Pi} \left\{ \hat{Y}_{(2)DE} \right\} = Y$ , i.e., the estimator  $\hat{Y}_{(2)R}$  is negatively biased and must be increased. To this purpose the  $x_{jS}$  are likely to provide an auxiliary information

well correlated with the survey variable. Therefore, the weights in (SI.17) can be calibrated in such a way that

$$\sum_{j \in \Pi_R} \gamma w_j x_j = X \quad (\text{SI.18})$$

i.e., the calibrated estimator of  $X$  performed on the respondent sample invariably guesses the true HSS total. Solving the calibration equation (SI.18) with respect to  $\gamma$  we achieve

$$\gamma = \frac{X}{\sum_{j \in \Pi_R} w_j x_j} = \frac{X}{\hat{X}_{(2)R}}$$

where, in analogy with (SI.17),  $\hat{X}_{(2)R}$  is the estimator of  $X$  based on the respondent sample  $\Pi_R$  and as such it is also negatively biased. Therefore, using the new calibrated weights  $\gamma w_j$  into (SI.17), the calibrated estimator turns out to be

$$\hat{Y}_{(2)DE-CAL} = \hat{Y}_{(2)R} \frac{X}{\hat{X}_{(2)R}} \quad (\text{SI.19})$$

As it is apparent from (SI.19),  $\hat{Y}_{(2)DE-CAL}$  is likely to provide an increase of the negatively biased estimator  $\hat{Y}_{(2)R}$ . In practice, the estimator arising from the nonresponse calibration weighting is analogous to a ratio estimator where the two estimators involved are performed on the respondent sample  $\Pi_R$ .

Being the ratio of two estimators, the calibrated estimator (SI.19) has an intractable mean and variance. However, arguing as in Fattorini et al (2013), it is apparent that  $\hat{Y}_{(2)R}$  and  $\hat{X}_{(2)R}$  are the unbiased DE estimators, analogous to (SI.9), of the totals of the variables  $z_j y_j$  and  $z_j x_j$  in the whole population  $U$ . In turn, these totals coincide with the totals  $Y_R$  and  $X_R$  of the  $y_j$ s and  $x_j$ s in the population of respondent cells  $U_R$ . Accordingly, following the Taylor series approximation up to the first order as in Särndal et al. (1992, Section 5.5), the estimator (SI.19) has an approximate expectation

$$E_{\Omega \circ \Theta \circ \Pi} \left\{ \hat{Y}_{(2)DE-CAL} \right\} \sim Y_R \frac{X}{X_R} \quad (\text{SI.20})$$

After some straightforward computations, it follows that the approximate expectation (SI.20) is equal to  $Y$  when the ratio of the two totals in the respondent cells  $\beta_R = Y_R/X_R$  is approximately equal to the ratio of the two totals in the population of nonrespondent cells. Hence, in this case, the calibration estimation (SI.19) is approximately unbiased with approximate variance

$$Var_{\Omega \circ \Theta \circ \Pi} \left[ \hat{Y}_{(2)DE-CAL} \right] \sim \left( \frac{X}{X_R} \right)^2 Var_{\Omega \circ \Theta \circ \Pi} \left[ \hat{E}_{(2)R} \right] \quad (\text{SI.21})$$

where  $\hat{E}_{(2)R}$  is the DE estimator of type (SI.9) of the total of the  $e_j$ s, where  $e_j = z_j (y | j - \beta_R x_j)$ , while, in analogy with (SI.11),

$$Var_{\Omega \circ \Theta \circ \Pi} \left[ \hat{E}_{(2)R} \right] = \sum_{l=1}^m \left[ X(Q_l) \sum_{q \in Q_l} \frac{E^2(q)}{X(q)} - E^2(Q_l) \right]$$

$$+ E_{\Omega} \left[ \sum_{l=1}^m \frac{X^2(Q_l)}{X^2(q_l)} \text{Var}_{\Theta} \left\{ \hat{E}_{(1,2)}(q_l) \mid \Omega \right\} \right] + E_{\Omega \circ \Theta} \left[ \sum_{l=1}^m \frac{X^2(Q_l)}{X^2(q_l)} \text{Var}_{\Pi} \left\{ \hat{E}_{(2)}(q_l) \mid \Omega, \Theta \right\} \right]$$

where  $E(Q_l)$ ,  $E(q_l)$  and  $E(q_l, i)$  are the totals of the  $e_j$ s within the q-block  $Q_l$ , the quadrat  $q_l$  and the  $i$ -th c-block of  $q_l$ , respectively,  $\hat{E}_{(1,2)}(q_l)$  is the second-stage estimator of  $E(q_l)$  and

$$\text{Var}_{\Theta} \left\{ \hat{E}_{(1,2)}(q_l) \mid \Omega \right\} = \sum_{i \in C(q_l)} X(q_l, i) \sum_{j \in U(q_l, i)} \frac{e_j}{x_j} - \sum_{i \in C(q_l)} E^2(q_l, i)$$

while

$$\begin{aligned} & \text{Var}_{\Pi} \left[ \hat{E}_{(2)}(q_l) \mid \Omega, \Theta \right] \\ &= \sum_{h > i \in C(q_l)} \left( \pi_{j(q_l, i)} \mid \pi_{j(q_l, h)} - \pi_{j(q_l, i), j(q_l, h)} \right) \left[ \frac{X(q_l, i) e_{j(q_l, i)}}{x_{j(q_l, i)} \pi_{j(q_l, i)}} - \frac{X(q_l, h) e_{j(q_l, h)}}{x_{j(q_l, h)} \pi_{j(q_l, h)}} \right]^2 \end{aligned}$$

Then, in analogy with equation (SI.12), the HH-like variance estimator is given by

$$V_{DE-CAL}^2 = \frac{1}{n(n-1)} \sum_{l=1}^m \sum_{i \in \Pi(q_l)} t_{j(q_l, i)}^2 \quad (\text{SI.22})$$

where, in this case,

$$t_{j(q_l, i)} = n \frac{\hat{e}_{j(q_l, i)}}{\tau_{j(q_l, i)}}$$

with  $\hat{e}_j = z_j (y_j - \hat{\beta}_R x_j)$  and  $\hat{\beta}_R = \hat{Y}_{(2)R} / \hat{X}_{(2)R}$ .

To exploit the auxiliary information achieved from previous investigation we have also attempted to exploit the HSSs for calibrating, mutatis mutandis, the difference estimator, i.e.,

$$\hat{Y}_{(2)DIF-CAL} = H + \hat{D}_{(2)R} \frac{X}{\hat{X}_{(2)R}} \quad (\text{SI.23})$$

where  $\hat{D}_{(2)R}$  is the unbiased DE estimator, analogous to (SI.9), of the total of the variables  $z_j d_j$  that coincides with the total  $D_R$  of the  $d_j$ s in the population of respondent cells  $U_R$ . Therefore, in analogy with the estimator  $\hat{Y}_{(2)DE-CAL}$ , the estimator  $\hat{Y}_{(2)DIF-CAL}$  is approximately unbiased if the ratio  $\delta_R = D_R / X_R$  is approximately equal to the ratio of the two totals in the population of nonrespondent cells. Moreover, the approximate variance is given by equation (SI.21) with  $e_j = z_j (d_j - \delta_R x_j)$ , while the HH-like variance estimator is given by

$$V_{DIF-CAL}^2 = \frac{1}{n(n-1)} \sum_{l=1}^m \sum_{i \in \Pi(q_l)} t_{j(q_l, i)}^2 \quad (\text{SI.24})$$

where, in this case,

$$t_{j(q_i,i)} = n \frac{\hat{e}_{j(q_i,i)}}{\tau_{j(q_i,i)}}$$

with  $\hat{e}_j = z_j \left( d_j - \hat{\delta}_R x_j \right)$  and  $\hat{\delta}_R = \hat{D}_{(2)R} / \hat{X}_{(2)R}$ .

## 8 Estimation within domains

It may occur that coverage estimates of a habitat are required not only for the whole study region but also for  $L$  portions partitioning the study region, usually referred to as domains in the sampling literature (e.g., Särndal et al., 1992, Chapter 10). For example, the study region may be divided into biogeographical regions, or by administrative regions and we could be interested to estimate the habitat coverage for each of the  $L$  regions.

Denote by  $U_1, \dots, U_L$  the  $L$  subpopulations of cells partitioning  $U$ , whose coverages are of interest. Following Särndal et al. (1992, section 10.3), estimation within domains can be performed simply by introducing for each domain  $g = 1, \dots, L$ , the domain indicator  $u_{g,j}$  equal to 1 for each cell  $j \in U_g$  and equal to 0 otherwise and then by applying the complete sample DE and DIF estimators (SI.9) and (SI.13) or the calibration weighting versions (SI.19) and (SI.23) and the corresponding variance estimators to the sample values involved multiplied by the domain indicator  $u_{g,j}$ . As the domain becomes smaller and smaller, the sample information available to make estimation reduces, deteriorating the precision of the domain estimators. The issue is well known in sampling literature as the small area estimation problem (e.g., Rao and Molina 2015, and references therein).

## References

- Brewer KRW, Hanif M (1983) Sampling with Unequal Probabilities. Springer, New York.
- Chen XH, Dempster AP, Liu JS (1994) Weighted finite population sampling to maximize entropy. *Biometrika* 81, 457–469.
- Deville JC, Tillé Y (1998) Unequal probability sampling without replacement through a splitting method. *Biometrika* 85, 89–101.
- Fattorini L, Franceschi S, Maffei D (2013) Design-based treatment of unit nonresponse in environmental surveys using calibration weighting. *Biometrical Journal* 55, 925–943.
- Haziza D, Thompson KJ, Yung W (2010) The effect of nonresponse adjustments on variance estimation. *Survey Methodology* 36, 35–43.
- Hedayat AS, Sinha BK (1991) Design and Inference in Finite Population Sampling. Wiley, New York.
- Little RJA, Rubin DB (2002) Statistical Analysis with Missing Data (2nd ed). Wiley, New York.
- Rao JNK, Molina I (2015) Small Area Estimation. Wiley, New York.
- Sampford MR (1967) On sampling without replacement with unequal probabilities of selection. *Biometrika* 54, 499–513.
- Särndal CE, Lundström S (2005) Estimation in Survey with Nonresponse. Wiley, New York.
- Särndal CE, Swensson B, Wretman J (1992) Model Assisted Survey Sampling. Springer, Berlin.

Tillé Y (1996) An elimination procedure for unequal probability sampling without replacement. *Biometrika* 83, 238–241.

Wolter KM (2007). *Introduction to Variance Estimation* (2nd ed). Springer, New York