

A comparison of design-based and model-based approaches for finite population spatial data.

Michael Dumelle^{*,a}, Matt Higham^b, Jay M. Ver Hoef^c, Anthony R. Olsen^a, Lisa Madsen^d

^a*United States Environmental Protection Agency, 200 SW 35th St, Corvallis, Oregon, 97333*

^b*Saint Lawrence University Department of Mathematics, Computer Science, and Statistics, 23 Romoda Drive, Canton, New York, 13617*

^c*Marine Mammal Laboratory, Alaska Fisheries Science Center, National Oceanic and Atmospheric Administration, Seattle, Washington, 98115*

^d*Oregon State University Department of Statistics, 239 Weniger Hall, Corvallis, Oregon, 97331*

Abstract

1. The design-based and model-based approaches to frequentist statistical inference rest on fundamentally different foundations. In the design-based approach, inference relies on random sampling. In the model-based approach, inference relies on distributional assumptions. We compare the approaches for finite population spatial data.
2. We provide relevant background for the design-based and model-based approaches and then study their performance using simulated and real data. In the simulated and real data, a variety of sample sizes, location layouts, dependence structures, and response types are considered. The population mean is the parameter of interest and performance is measured using statistics like bias, squared error, and interval coverage.
3. When studying the simulated and real data, we found that regardless of the strength of spatial dependence in the data, the generalized random tessellation stratified (GRTS) algorithm, which explicitly incorporates spatial locations into sampling, tends to outperform the simple random sampling (SRS) algorithm, which does not explicitly incorporate spatial locations into sampling. We also found that model-based approaches tend

^{*}Corresponding Author: Michael Dumelle (Dumelle.Michael@epa.gov)

Preprint submitted to *Methods in Ecology and Evolution*

April 13, 2022

to outperform design-based approaches, even for skewed data where the model-based distributional assumptions are violated. The performance gap between these approaches is small GRTS samples are used but large when SRS samples are used. This suggests that the sampling choice (whether to use GRTS or SRS) is most important when performing design-based inference.

4. There are many benefits and drawbacks to the design-based and model-based approaches for finite population spatial data that practitioners must consider when choosing between them. We provide relevant background contextualizing each approach and study their properties in a variety of scenarios, making recommendations for use based on the practitioner's goals.

Keywords

Design-based inference; Finite population block kriging (FPBK); Generalized random tessellation stratified (GRTS) algorithm; Local neighborhood variance estimator; Model-based inference; Restricted maximum likelihood (REML) estimation; Spatially balanced sampling; Spatial covariance

1. Introduction

When data cannot be collected for all units in a population (i.e., population units), data are collected on a subset of the population units – this subset is called a sample. There are two general approaches for using samples to make frequentist statistical inferences about a population: design-based and model-based. In the design-based approach, inference relies on randomly assigning some population units to be in the sample (random sampling). Alternatively, in the model-based approach, inference relies on distributional assumptions about

55 the underlying stochastic process generating the sample. Each paradigm has a
 56 deep historical context (Sterba, 2009) and its own set of benefits and drawbacks
 57 (Hansen et al., 1983, p. @brus1997random). In this manuscript, we compare the
 58 design-based and model-based approaches for finite population spatial data.

59 Spatial data are data that have some sort of spatial index, usually via
 60 coordinates. De Gruijter and Ter Braak (1990) and Brus and DeGruijter (1993)
 61 give early comparisons of design-based and model-based approaches for spatial
 62 data, quashing the belief that design-based approaches could not be used for
 63 spatially correlated data. Since then, there have been several general comparisons
 64 between design-based and model-based approaches for spatial data (Brus and
 65 De Gruijter, 1997; Brus, 2021; Ver Hoef, 2002, 2008). Cooper (2006) reviews the
 66 two approaches in an ecological context before introducing a “model-assisted”
 67 variance estimator that combines aspects from each approach. In addition
 68 to Cooper (2006), there has been substantial research and development into
 69 estimators that use both design-based and model-based principles (see e.g., Sterba
 70 (2009) and Cicchitelli and Montanari (2012), and for Bayesian approaches, see
 71 Chan-Golston et al. (2020) and Hofman and Brus (2021)).

72 While comparisons between design-based and model-based approaches have
 73 been studied in spatial contexts, our contribution is comparing design-based
 74 approaches specifically built for spatial data to model-based approaches. Though
 75 the broad comparisons we draw between design-based and model-based ap-
 76 proaches generalize to finite and infinite populations, we focus on finite popu-
 77 lations. A finite population contains a finite number of population units (we
 78 assume the finite number is known); an example is lakes (treated as a whole with
 79 the lake centroid representing location) in the conterminous United States. An
 80 infinite population contains an infinite number of population units; an example
 81 is locations within a single lake.

82 The rest of the manuscript is organized as follows. In Section 1.1, we introduce
 83 and provide relevant background for design-based and model-based approaches
 84 to finite population spatial data. In Section 2, we describe how we intend to
 85 compare performance of the approaches using simulated and real data. In Section
 86 3, we present analysis results for the simulated and real data. And in Section 4,
 87 we end with a discussion and provide directions for future research.

88 *1.1. Background*

89 The design-based and model-based approaches incorporate randomness in
 90 fundamentally different ways. In this section, we describe the role of randomness
 91 for each approach and the subsequent effects on statistical inferences for spatial
 92 data.

93 *1.1.1. Comparing Design-Based and Model-Based Approaches*

94 The design-based approach assumes the population is fixed. Randomness
 95 is incorporated via the selection of population units according to a sampling
 96 design. A sampling design assigns a probability of selection to each sample
 97 (subset of population units). Some examples of commonly used sampling designs
 98 include simple random sampling, stratified random sampling, and cluster sam-
 99 pling. The inclusion probability of a population unit follows by summing each
 100 sample's probability of selection over all samples that contain the population
 101 unit. Inclusion probabilities are later used to estimate population parameters.

102 When samples are chosen in a manner such that the layout of sampled units
 103 reflects the layout of the population units, we call the resulting sample spatially
 104 balanced. By “reflecting the layout of the population units”, we mean that if
 105 population units are concentrated in specific areas, the units in the sample should
 106 be concentrated in the same areas. Because spatially balanced samples reflect
 107 the layout of the population units, they are not necessarily spread out in space in

108 some equidistant manner. One approach to selecting spatially balanced samples
 109 is the generalized random tessellation stratified (GRTS) algorithm (Stevens and
 110 Olsen, 2004), which we discuss in more detail in Section 1.1.2.

111 Fundamentally, the design-based approach combines the randomness of the
 112 sampling design with the data collected via the sample to justify the estimation
 113 and uncertainty quantification of fixed, unknown parameters of a population (e.g.,
 114 a population mean). Treating the data as fixed and incorporating randomness
 115 through the sampling design yields estimators having very few other assumptions.
 116 Confidence intervals for these types of estimators are typically derived using
 117 limiting arguments that incorporate all possible samples. Sample means, for
 118 example, are asymptotically normal (Gaussian) by the Central Limit Theorem
 119 (under some assumptions). If we repeatedly select samples from the population,
 120 then 95% of all 95% confidence intervals constructed from a procedure with
 121 appropriate coverage will contain the true fixed population mean. Särndal et al.
 122 (2003) and Lohr (2009) provide thorough reviews of the design-based approach.

123 The model-based approach assumes the population is a random realization of a
 124 data-generating stochastic process. Randomness is formally incorporated through
 125 distributional assumptions on this process. Strictly speaking, randomness need
 126 not be incorporated through random sampling, though Diggle et al. (2010) warn
 127 against preferential sampling. Preferential sampling occurs when the process
 128 generating the data locations and the process being modeled are not independent
 129 of one another. To guard against preferential sampling, model-based approaches
 130 can implement some form of random sampling, though it is common for model-
 131 based approaches to sample non-randomly. When model-based approaches
 132 do implement random sampling, the inclusion probabilities are ignored when
 133 analyzing the sample (in contrast to the design-based approach, which relies on
 134 these inclusion probabilities to analyze the sample).

135 Instead of estimating fixed, unknown population parameters, as in the design-
136 based approach, often the goal of model-based inference is to predict a realized
137 variable. For example, suppose the realized mean of all population units (the
138 realized population mean) is the variable of interest. Instead of a fixed, unknown
139 mean, we are predicting the value of the mean, a random variable. Prediction
140 intervals are then derived using assumptions of the data-generating stochastic
141 process. If we repeatedly generate realizations from the same process and select
142 samples, then 95% of all 95% prediction intervals constructed from a procedure
143 with appropriate coverage will contain their respective realized means. Cressie
144 (1993) and Schabenberger and Gotway (2017) provide thorough reviews of model-
145 based approaches for spatial data. In Fig. 1, we provide a visual comparison
146 of the design-based and model-based approaches (Ver Hoef (2002) and Brus
147 (2021) provide similar figures). This figure contrasts the design-based approach
148 with a fixed population and random sampling to the model-based approach with
149 random populations and non-random sampling.

150 1.1.2. *Spatially Balanced Design and Analysis*

151 We previously mentioned that the design-based approach can be used to
152 select spatially balanced samples. Spatially balanced samples are useful because
153 parameter estimates from these samples tend to vary less than parameter es-
154 timates from samples lacking spatial balance (Barabesi and Franceschi, 2011;
155 Benedetti et al., 2017; Grafström and Lundström, 2013; Robertson et al., 2013;
156 Stevens and Olsen, 2004; Wang et al., 2013). The first spatially balanced sam-
157 pling algorithm to see widespread use was the generalized random tessellation
158 stratified (GRTS) algorithm (Stevens and Olsen, 2004). To quantify the spatial
159 balance of a sample, Stevens and Olsen (2004) proposed loss metrics based on
160 Voronoi polygons (i.e., Dirichlet Tessellations). After the GRTS algorithm was
161 developed, several other spatially balanced sampling algorithms emerged, includ-

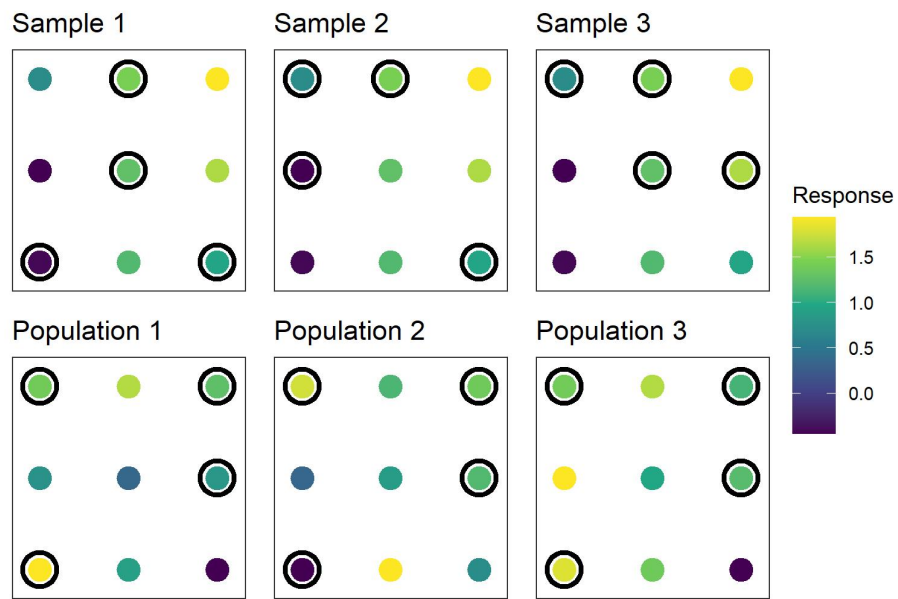


Figure 1: A visual comparison of the design-based and model-based approaches. In the top row, the design-based approach is highlighted. There is one fixed population with nine population units and three random samples of size four (points circled are those sampled). The response values at each site are fixed. In the bottom row, the model-based approach is highlighted. There are three realizations of the same data-generating stochastic process that are all sampled at the same four locations. The response values at each site are random.

162 ing stratified sampling with compact geographical strata Walvoort et al. (2010),
 163 the local pivotal method (Grafström et al., 2012; Grafström and Matei, 2018),
 164 spatially correlated Poisson sampling (Grafström, 2012), balanced acceptance
 165 sampling (Robertson et al., 2013), within-sample-distance sampling (Benedetti
 166 and Piersimoni, 2017), and Halton iterative partitioning sampling (Robertson
 167 et al., 2018). In this manuscript, we select spatially balanced samples using
 168 the GRTS algorithm because it is readily available in the **spsurvey R** package
 169 (Dumelle et al., 2022) and naturally accommodates finite and infinite sampling
 170 frames, unequal inclusion probabilities, and replacement units. Replacement
 171 units are additional population units that can be sampled when a population unit
 172 originally selected can no longer be sampled. A couple reasons why an originally
 173 selected site can no longer be sampled include its location being physically
 174 inaccessible or on private land that the researcher does not have permission to
 175 access.

176 The GRTS algorithm selects samples by utilizing a particular mapping
 177 between two-dimensional and one-dimensional space that preserves proximity
 178 relationships. First the bounding box of the domain is split up into four distinct,
 179 equally sized squares called level-one cells. Each level-one cell is randomly
 180 assigned a level-one address of 0, 1, 2, or 3. The set of level-one cells is denoted
 181 by \mathcal{A}_1 and defined as $\mathcal{A}_1 \equiv \{a_1 : a_1 = 0, 1, 2, 3\}$. Within each level-one cell, the
 182 inclusion probability for each population unit is summed, and if any of these
 183 sums exceed one, a second level of cells is added. Then each level-one cell is split
 184 into four distinct, equally sized squares called level-two cells. Each level-two cell
 185 is randomly assigned a level-two address of 0, 1, 2, or 3. The set of level-two
 186 cells is denoted by \mathcal{A}_2 and defined as $\mathcal{A}_2 \equiv \{a_1 a_2 : a_1 = 0, 1, 2, 3; a_2 = 0, 1, 2, 3\}$.
 187 The inclusion probabilities within each level-two cell are summed, and if any of
 188 these sums exceed one, a third level of cells is added. This process continues for

189 k steps, until all level- k cells have inclusion probability sums no larger than one.

190 Then $\mathcal{A}_k \equiv \{a_1 \dots a_k : a_1 = 0, 1, 2, 3; \dots; a_k = 0, 1, 2, 3\}$.

191 After determining \mathcal{A}_k , it is placed into hierarchical order. Hierarchical order
 192 is a numeric order that first sorts \mathcal{A}_k by the level-one addresses from smallest
 193 to largest, then sorts \mathcal{A}_k by the level-two addresses from smallest to largest, and so
 194 on. For example, \mathcal{A}_2 in hierarchical order is the set
 195 $\{00, 01, 02, 03, 10, \dots, 13, 20, \dots, 23, 30, \dots, 33\}$. Because hierarchical ordering sorts
 196 by level-one cells, then level-two cells, and so on, population units that have
 197 similar hierarchical addresses tend to be nearby one another in space. Next each
 198 population unit is mapped to a one-dimensional line in hierarchical order where
 199 each population unit's inclusion probability equals its line-length. If a level- k
 200 cell has multiple population units in it, they are randomly placed within the
 201 cell's respective line segment. A uniform random variable is then simulated in
 202 $[0, 1]$ and a systematic sample is selected on the line, yielding n sample points for
 203 a sample size n . Each of these sample points falls on some population unit's line
 204 segment, and thus that population unit is selected in the sample. For further
 205 details regarding the GRTS algorithm, see Stevens and Olsen (2004).

After selecting a sample and collecting data, unbiased estimates of population means and totals can be obtained using the Horvitz-Thompson estimator (Horvitz and Thompson, 1952). If τ is a population total, the Horvitz-Thompson estimator for τ , denoted by $\hat{\tau}_{ht}$, is given by

$$\hat{\tau}_{ht} = \sum_{i=1}^n z_i \pi_i^{-1}, \quad (1)$$

206 where z_i is the value of the i th population unit in the sample, π_i is the inclusion
 207 probability of the i th population unit in the sample, and n is the sample size. An
 208 estimate of the population mean is obtained by dividing $\hat{\tau}_{ht}$ by N , the number
 209 of population units.

210 It is also important to quantify the uncertainty in $\hat{\tau}_{ht}$. Horvitz and Thompson
 211 (1952) and Sen (1953) provide variance estimators for $\hat{\tau}_{ht}$, but these estimators
 212 have two drawbacks. First, they rely on calculating π_{ij} , the probability that
 213 population unit i and population unit j are both in the sample – this quantity
 214 can be challenging if not impossible to calculate analytically for GRTS samples.
 215 Second, these estimators tend to ignore the spatial locations of the population
 216 units. To address these two drawbacks simultaneously, Stevens and Olsen (2003)
 217 proposed the local neighborhood variance estimator. The local neighborhood
 218 variance estimator does not rely on π_{ij} and estimates the variance of $\hat{\tau}$ conditional
 219 on the random properties of the GRTS sample – the idea being that this
 220 conditioning should yield a more precise estimate of $\hat{\tau}$. They show that the
 221 contribution from each sample unit (population unit in the sample) to the overall
 222 variance is dominated by local variation. Thus the local neighborhood variance
 223 estimator is a weighted sum of variance estimates from each sample unit’s local
 224 neighborhood. These local neighborhoods contain the sample unit itself and
 225 its three nearest neighbors among all other sample units. For more details, see
 226 Stevens and Olsen (2003).

227 1.1.3. Finite Population Block Kriging

228 Finite population block kriging (FPBK) is a model-based approach that
 229 expands the geostatistical Kriging framework to the finite population setting
 230 (Ver Hoef, 2008). Instead of developing inference based on a specific sampling
 231 design, we assume the data are generated by a spatial stochastic process. We
 232 summarize some of the basic principles of FPBK next – for more details, see
 233 Ver Hoef (2008). Let $\mathbf{z} \equiv \{z(s_1), z(s_2), \dots, z(s_N)\}$ be an $N \times 1$ response vector
 234 at locations s_1, s_2, \dots, s_N that can be measured at the N population units.
 235 Suppose we want to use a sample to predict some linear function of the response
 236 variable, $f(\mathbf{z}) = \mathbf{b}'\mathbf{z}$, where \mathbf{b}' is a $1 \times N$ vector of weights (e.g, the population

mean is represented by a weights vector whose elements all equal $1/N$). Denoting quantities that are part of the sampled population units with a subscript s and quantities that are part of the unsampled population units with a subscript u , let

$$\begin{pmatrix} \mathbf{z}_s \\ \mathbf{z}_u \end{pmatrix} = \begin{pmatrix} \mathbf{X}_s \\ \mathbf{X}_u \end{pmatrix} \boldsymbol{\beta} + \begin{pmatrix} \boldsymbol{\delta}_s \\ \boldsymbol{\delta}_u \end{pmatrix}, \quad (2)$$

where \mathbf{X}_s and \mathbf{X}_u are the design matrices for the sampled and unsampled population units, respectively, $\boldsymbol{\beta}$ is the parameter vector of fixed effects, and $\boldsymbol{\delta} \equiv [\boldsymbol{\delta}_s \ \boldsymbol{\delta}_u]'$, where $\boldsymbol{\delta}_s$ and $\boldsymbol{\delta}_u$ are random errors for the sampled and unsampled population units, respectively.

FPBK assumes $\boldsymbol{\delta}$ in Equation 2 has mean-zero and a spatial dependence structure that can be modeled using a covariance function. This covariance function is commonly assumed to be non-negative, second-order stationary (depending only on the separation vector (e.g., distance) between population units), isotropic (independent of direction), and decay with distance between population units (Cressie, 1993). Henceforth, it is implied that we have made these same assumptions regarding $\boldsymbol{\delta}$, though Chiles and Delfiner (1999), pp. 80-93 discuss covariance functions that are not second-order stationary, not isotropic, or not either. A variety of flexible covariance functions can be used to model $\boldsymbol{\delta}$ (Cressie, 1993); one example is the exponential covariance function (Cressie (1993) provides a thorough list of spatial covariance functions). The i, j th element of the exponential covariance matrix, $\text{cov}(\boldsymbol{\delta})$, is

$$\text{cov}(\delta_i, \delta_j) = \begin{cases} \sigma_1^2 \exp(-h_{i,j}/\phi) & h_{i,j} > 0 \\ \sigma_1^2 + \sigma_2^2 & h_{i,j} = 0 \end{cases}, \quad (3)$$

where σ_1^2 is the variance parameter that quantifies the spatially dependent

246 variability, σ_2^2 is the variance parameter that quantifies that spatially independent
 247 variability, ϕ is the distance parameter that measures the distance-decay rate of
 248 the covariance, and $h_{i,j}$ is the Euclidean distance between population units i
 249 and j . In geostatistical literature, σ_1^2 is often called the partial sill, σ_2^2 is often
 250 called the nugget, and ϕ is often called the range.

The parameters in Equation 2 can be estimated using a variety of techniques,
 but we focus on using restricted maximum likelihood (Harville, 1977; Patterson
 and Thompson, 1971; Wolfinger et al., 1994). REML is preferred over maximum
 likelihood (ML) because ML estimates can be badly biased for small sample sizes,
 due to the fact that ML makes no adjustment for the simultaneous estimation of
 β and δ (Patterson and Thompson, 1971). Minus twice the REML log-likelihood
 of the sampled sites is given by

$$\ln |\Sigma| + (z_s - X_s \tilde{\beta})^T \Sigma_{ss}^{-1} (z_s - X_s \tilde{\beta}) + \ln |X_s^T \Sigma_{ss}^{-1} X_s| + (n - p) \ln(2\pi), \quad (4)$$

251 where $\tilde{\beta} = (X_s^T \Sigma_{ss}^{-1} X_s)^{-1} X_s^T \Sigma_{ss}^{-1} z_s$ and Σ_{ss} is the covariance matrix of the
 252 sampled sites. Minimizing Equation 4 yields $\hat{\delta}_{reml}$, the REML estimates of
 253 δ . Then β_{reml} , the REML estimate of β , is given by $(X_s^T \hat{\Sigma}_{ss}^{-1} X_s)^{-1} X_s^T \hat{\Sigma}_{ss}^{-1} z_s$,
 254 where $\hat{\Sigma}_{ss}$ is Σ_{ss} evaluated at $\hat{\delta}_{reml}$.

255 With the model formulation in Equation 2, the Best Linear Unbiased Predictor
 256 (BLUP) for $f(\mathbf{b}'\mathbf{z})$ and its prediction variance can be computed. While details
 257 of the derivation are in Ver Hoef (2008), we note here that the predictor and
 258 its variance are both moment-based, meaning that they do not rely on any
 259 distributional assumptions. Distributional assumptions are used, however, when
 260 constructing prediction intervals.

261 Other approaches, such as k-nearest-neighbors (Fix and Hodges, 1989; Ver
 262 Hoef and Temesgen, 2013) and random forest (Breiman, 2001), among others,
 263 could also be used to obtain predictions for a mean or total from finite population

264 spatial data. Compared to the k-nearest-neighbors and random forest approach,
 265 we prefer FPBK because it is model-based and relies on theoretically-based
 266 variance estimators leveraging the model's spatial covariance structure, whereas
 267 k-nearest-neighbors and random forests use ad-hoc variance estimators (Ver Hoef
 268 and Temesgen, 2013). Additionally, Ver Hoef and Temesgen (2013) compared
 269 FPBK, k-nearest-neighbors, and random forest in a variety of spatial data
 270 contexts, and FPBK tended to perform best.

271 2. Materials and Methods

In this section we describe how we used simulated and real data to investigate performance between simple random sampling without replacement (SRS) and GRTS sampling as well as performance between design-based (DB) and model-based (MB) inference. In SRS and GRTS sampling, all population units had equal inclusion probabilities. The important distinction between SRS and GRTS is that SRS ignores spatial locations while sampling but GRTS explicitly incorporates them. Together, the two sampling plans (SRS and GRTS) combined with the two inference approaches (DB and MB) yielded four sampling-inference combinations: SRS-DB, SRS-MB, GRTS-DB, and GRTS-MB. For SRS-DB, the Horvitz-Thompson estimator (1) was used to estimate means and the commonly-used SRS variance formula (Lohr, 2009; Särndal et al., 2003) was used to estimate the variance. This variance formula is given by

$$\frac{f[\sum_{i=1}^n (z_i - \bar{z})^2]}{n(n-1)}, \quad (5)$$

272 where z_i is the i th response value, \bar{z} is the mean of all z_i , n is the sample size, N
 273 is the population size, and $f = (1 - n/N)$ (f is often called the finite population
 274 correction factor). For GRTS-DB, the Horvitz-Thompson estimator was used
 275 to estimate means and the local neighborhood variance was used to estimate

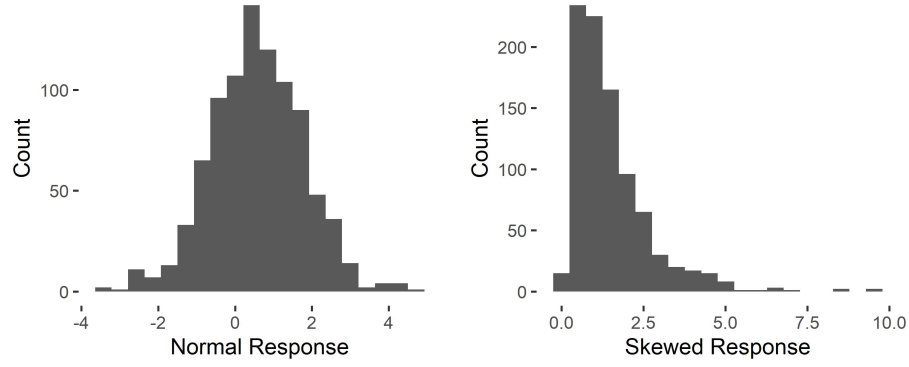
variances. For SRS-MB and GRTS-MB, FPBK was used to estimate means and variances and parameters were estimated using restricted maximum likelihood.

We used simulated data to compare the sampling-inference combinations across many realized populations from the same data-generating stochastic process. With the simulated data, we were in control of the data-generating stochastic process and the random sampling process. We used real data from the 2012 National Lakes Assessment (USEPA, 2012) to compare the sampling-inference combinations within a single realized population (which is typically the case in reality). With the real data, we were in control of only the random sampling process.

2.1. Simulated Data

CHANGE LOGNORMAL VERBAGE TO SKEWED – look for DRE acronym

We evaluated performance of the four sampling-inference combinations in 36 different simulation scenarios. The 36 scenarios resulted from the crossing of three sample sizes, two location layouts (of the population units), two response types, and three proportions of dependent random error (DRE). The three sample sizes (n) were $n = 50$, $n = 100$, and $n = 200$. Samples were always selected from a population size (N) of $N = 900$. The two location layouts were random and gridded. Locations in the random layout were randomly generated inside the unit square $([0, 1] \times [0, 1])$. Locations in the gridded layout were placed on a fixed, equally spaced grid inside the unit square. The two response types were normal and skewed. For the normal response type, the response was simulated using mean-zero random errors with the exponential covariance (Equation 3) for three proportions of dependent random error (DRE): 0% DRE, 50% DRE, and 90% DRE. Recall the proportion of DRE is represented by $\sigma_1^2/(\sigma_1^2 + \sigma_2^2)$, where σ_1^2 and σ_2^2 are the DRE variance and independent random error (IRE) variance from Equation 3, respectively. The total variance, $\sigma_1^2 + \sigma_2^2$, was always 2. The distance



(a) Histogram of a realized population for the normal response. (b) Histogram of a realized population for the skewed response.

Figure 2: Histograms of realized populations simulated for the normal and skewed responses using the random layout and 50% DRE.

parameter was always $\sqrt{2}/3$, chosen so that the correlation in the DRE decayed to nearly zero at $\sqrt{2}$, the largest possible distance between two population units in the domain. For the skewed response type, the response was first simulated using the same approach as for the normal response type, except that the total variance was 0.6931 instead of 2. The response was then exponentiated, yielding a skewed random variable whose total variance was 2. The skewed responses were used to evaluate performance of the sampling-inference approaches for data that were not normal but were still estimated using REML, which relies on a normal log-likelihood. Figure 2 shows an example of a realized population for the normal and skewed responses using the random layout and 50% DRE.

In each of the 36 simulation scenarios, there were 2000 independent simulation trials. Within each simulation scenario and trial, IRS and GRTS samples were selected and then design-based and model-based analyses were used to estimate (design-based) or predict (model-based) the mean and construct 95% confidence (design-based) or 95% prediction (model-based) intervals. With the model-based analyses, covariance parameters were estimated (using REML) separately for each trial. After all 2000 trials, we summarized the long-run performance of

the sampling-inference combination in each scenario by calculating mean bias, root-mean-squared error, and interval coverage. Mean bias is taken as the average deviation between each trial's estimated (or predicted) mean and its realized mean: $\frac{1}{n} \sum_{i=1}^{2000} (\hat{\mu}_i - \mu_i)$, where i indexes simulation trials. Root-mean-squared error is taken as the square root of the average squared deviation between each trial's estimated (or predicted) mean and its realized mean: $\sqrt{\frac{1}{n} \sum_{i=1}^{2000} (\hat{\mu}_i - \mu_i)^2}$. Interval coverage is taken as the proportion of simulation trials where the realized mean is contained in its 95% confidence (or prediction) interval. These intervals are constructed using the normal distribution – justification comes from the asymptotic normality of means via the central limit theorem (under some assumptions). Quantifying these metrics is important because together, they give us an idea of the accuracy (mean bias), spread (RMSE), and validity (interval coverage) of the sampling-inference combinations.

2.2. National Lakes Assessment Data

The United States Environmental Protection Agency (USEPA), states, and tribes periodically conduct National Aquatic Research Surveys (NARS) to assess the water quality of various bodies of water in the contiguous United States. One component of NARS is the National Lakes Assessment (NLA), which measures various aspects of lake health and water quality. We will analyze mercury concentration data collected at 986 lakes from the 2012 NLA. Although we can calculate the true mean mercury concentration values for these 986 lakes, here we will explore whether or not we can obtain an adequately precise estimate (design-based) or prediction (model-based) for the realized mean mercury concentration if we sample only 100 of the 986 lakes. For each of the four familiar sampling-inference combinations (IRS-Design, IRS-Model, GRTS-Design, and GRTS-Model), we estimate (design-based) or predict (model-based) the mean mercury concentration and construct 95% intervals from this sample of 100 lakes and

compare to the true mean mercury concentration from all 986 lakes.

3. Results

3.1. Simulated Data

The mean bias was nearly zero for all four sampling-inference combinations in all 36 scenarios, so we omit a more detailed summary of those results here. Tables for mean bias in all 36 simulation scenarios are provided in the supporting information.

Fig. 3 shows the relative rMS(P)E of the four sampling analysis combinations using the random location layout with “IRS-Design” as the baseline. The relative rMS(P)E is defined as

$$\frac{\text{rMS(P)E of sampling-inference combination}}{\text{rMS(P)E of IRS-Design}},$$

When there is no spatial covariance (Fig. 3, “Prop DE: 0” row), the four sampling-inference combinations have approximately equal rMS(P)E. In these scenarios, using the GRTS algorithm or a model-based analysis does not increase efficiency compared to IRS-Design. When there is spatial covariance (Fig. 3, “Prop DE: 0.5” and “Prop DE: 0.9” rows), GRTS-Model tends to have the lowest rMS(P)E, followed by GRTS-Design, IRS-Model, and finally IRS-Design, though the difference in relative rMS(P)E among GRTS-Model, GRTS-Design, and IRS-Model is relatively small. As the strength of spatial covariance increases, the gap in rMS(P)E between IRS-Design and the other sampling-inference combinations widens. Finally we note that when there is spatial covariance, IRS-Model has a much lower rMS(P)E than IRS-Design, suggesting that the poor design properties of IRS are largely mitigated by the model-based analysis. These rMS(P)E conclusions are similar to those observed in the grid location

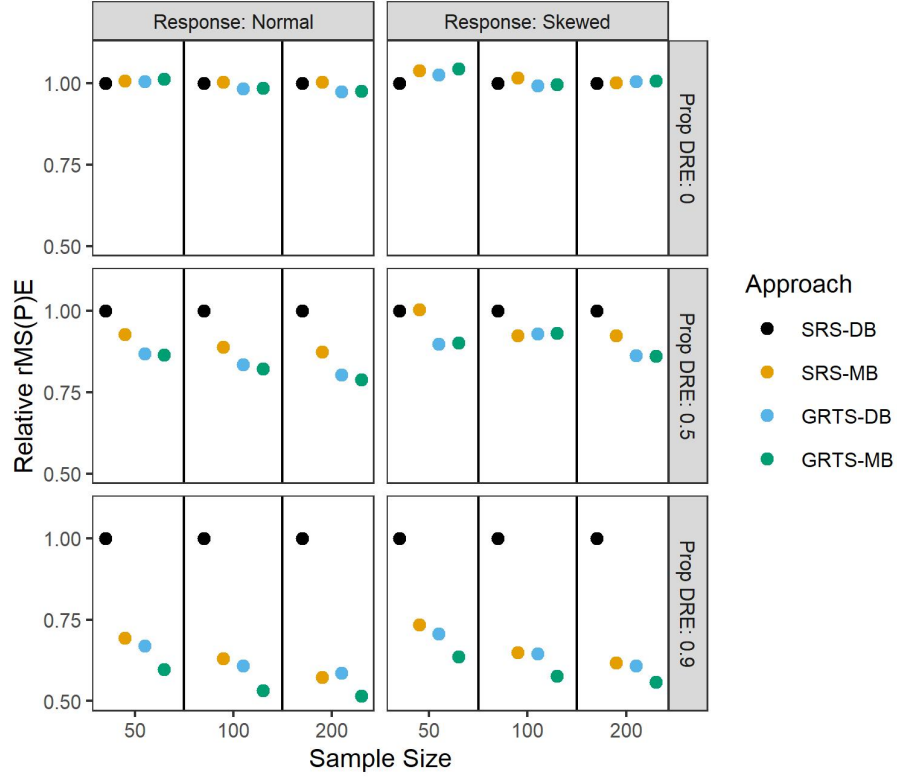


Figure 3: Relative rMS(P)E in the simulation study for the four sampling-inference combinations. The rows indicate the proportion of dependent error and the columns indicate the response type.

367 layout, so we omit a grid location layout figure here. Tables for rMS(P)E in all
 368 36 simulation scenarios are provided in the supporting information.

Fig. ?? shows the relative MStdE of the four sampling-inference combinations using the random location layout with “IRS-Design” as the baseline. The relative MStdE is defined as

$$\frac{\text{MStdE of sampling-inference combination}}{\text{MStdE of IRS-Design}},$$

369 Many general takeaways regarding MStdE are similar to general takeaways
 370 regarding rMS(P)E: there seems to be no benefit to using IRS, even when there

is no spatial covariance; as the strength of spatial covariance increases, the gap in MStdE between IRS-Design and the other sampling-inference combinations widens; and IRS-Model outperforms IRS-Design by a noticeable margin. These fact that the rMS(P)E and MStdE findings are similar is not particularly surprising because the mean bias for all sampling-inference combinations was nearly zero, thus rMS(P)E is driven by the standard error of the estimators (design-based) or predictors (model-based). We do note that between GRTS-Design and GRTS-Model, GRTS-Design had lower MStdE when there was no spatial covariance or a medium amount of spatial covariance (Fig. ??, “Prop DE: 0” and “Prop DE: 0.5” rows), and GRTS-Model had lower MStdE when there was a high amount of spatial covariance (Fig. ??, “Prop DE: 0.9” row). These MStdE conclusions are similar to those observed in the grid location layout, so we omit a grid location layout figure here. Tables for MStdE in all 36 simulation scenarios are provided in the supporting information.

Fig. 4 shows the 95% interval coverage for each of the four sampling-inference combinations in the random location layout. Within each scenario, the sampling-inference combinations tend to have fairly similar interval coverage, though when $n = 50$ or $n = 100$, GRTS-Design coverage is usually a few percentage points lower than the other combinations. Coverage in the normal response scenarios was usually near 95%, while coverage in the lognormal response scenarios usually varied from 90% to 95% but increased with the sample size. At a sample size of 200, all four sampling-inference combinations had approximately 95% interval coverage in both response scenarios for all dependent error proportions. These interval coverage conclusions are similar to those observed in the grid location layout, so we omit a grid location layout figure here. Tables for interval coverage in all 36 simulation scenarios are provided in the supporting information.

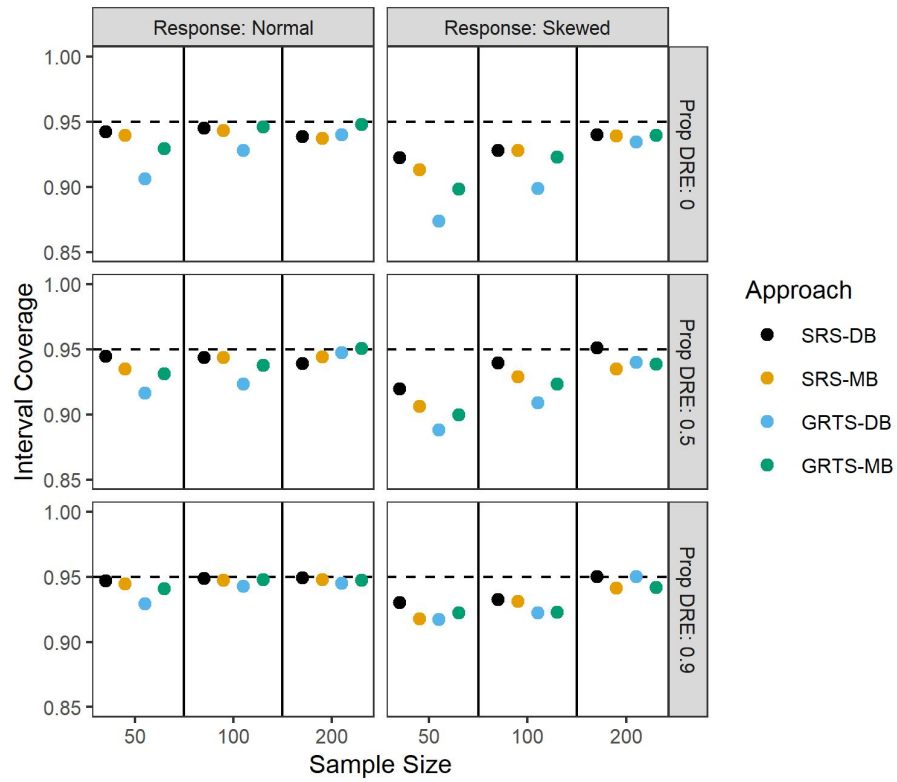


Figure 4: Interval coverage in the simulation study for the four sampling-inference combinations. The rows indicate the proportion of dependent error and the columns indicate the response type. The solid, black line represents 95% coverage.

3.2. National Lakes Assessment Data

USE MERCURY UNITS

Fig. ?? shows a map and histogram of mercury concentration in all 986 NLA lakes. The map shows mercury concentration exhibits some spatial patterning, with high mercury concentrations in the northeast and north central United States. The histogram shows that mercury concentration is right-skewed, with most lakes having a low value of mercury concentration but a few having a much higher concentration. Fig. ?? also shows mercury concentration's empirical semivariogram. The empirical semivariogram can be used as a tool to visualize spatial dependence. It quantifies the mean of the halved squared differences (semivariance) among all pairs of mercury concentrations at different distances apart. When a process has spatial covariance (exhibits spatial dependence), the mean semivariance tends to be smaller at small distances and larger at large distances. The empirical semivariogram in Fig. ?? suggests that mercury concentration exhibits spatial dependence. Lastly we note that the true mean mercury concentration in the 986 NLA lakes is 103.2 ng / g.

We selected a single IRS sample and a single GRTS sample and estimated (design-based) or predicted (model-based) the mean mercury concentration and constructed 95% confidence (design-based) and 95% (model-based) prediction intervals. For the model-based analyses, the exponential covariance was used. Table 1 shows the results from these analyses. Though we should not generalize these results to other samples from this population, we do mention a few findings. First, IRS-Design has the largest standard error. Second, compared to IRS-Design and IRS-Model, GRTS-Design and GRTS-Model are much closer to the true mean mercury concentration (have bias closer to zero) and have much lower standard errors (more precise intervals). Third, GRTS-Model has the least amount of bias and the lowest standard error (most precise interval). Finally,

we note that for all sampling-inference combinations, the true mean mercury concentration (103.2 ng / g) is within the bounds of the combination's 95% interval.

Approach	True Mean	Est/Pred	SE	95% LB	95% UB
IRS-Design	103.2	112.7	8.8	95.4	129.9
IRS-Model	103.2	110.5	7.9	95.0	125.9
GRTS-Design	103.2	101.8	6.1	89.8	113.7
GRTS-Model	103.2	102.3	5.9	90.8	113.9

Table 1: For each sampling-inference combination (Approach), the true mean mercury concentration (True Mean), estimates/predictions (Est/Pred), standard errors (SE), lower 95% interval bounds (95% LB), and upper 95% interval bounds (95% UB) for mean mercury concentration computed using a sample of 100 lakes in the NLA data.

3.3. New Application

4. Discussion

ADD EXTRAS LIKE ANISOTROPY AND UNEQUAL INCLUSION PROBABILITIES

The design-based and model-based approaches to statistical inference are fundamentally different paradigms. The design-based approach relies on random sampling to estimate population parameters. The model-based approach relies on distributional assumptions to predict realized values of a stochastic process. Though the model-based approach does not rely on random sampling, it can still be beneficial as a way to guard against preferential sampling. While the design-based and model-based approaches have often been compared in the literature from theoretical and analytical perspectives, our contribution lies in studying them in a spatial context while implementing spatially balanced sampling and the design-based, local neighborhood variance estimator. Aside from the theoretical differences described, a few analytical findings from the simulation study are particularly notable. First, independent of the analysis approach, we found no reason to prefer IRS over GRTS when sampling spatial

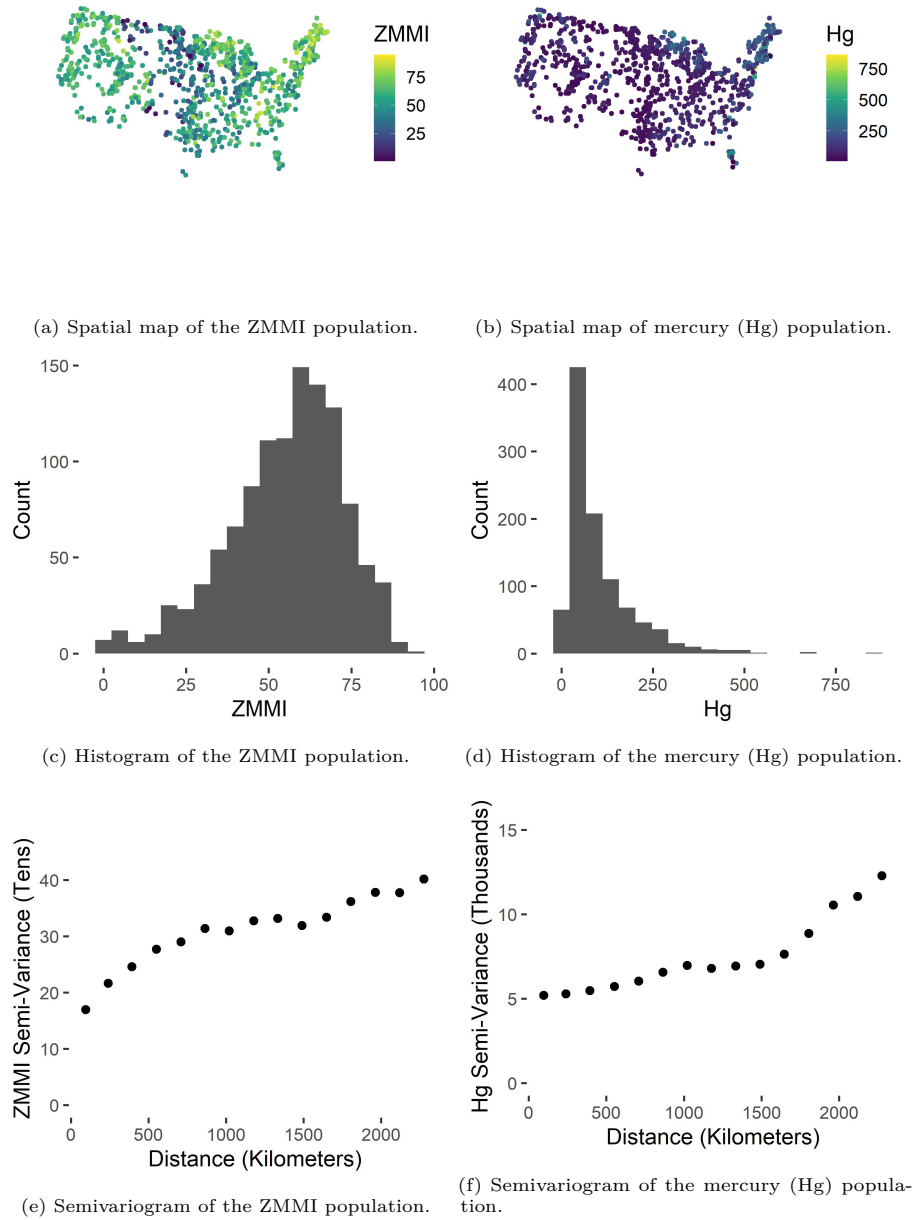


Figure 5: Exploratory graphics of the ZMMI and mercury (Hg) populations in the National Lakes Assessment (NLA) 2012 data.

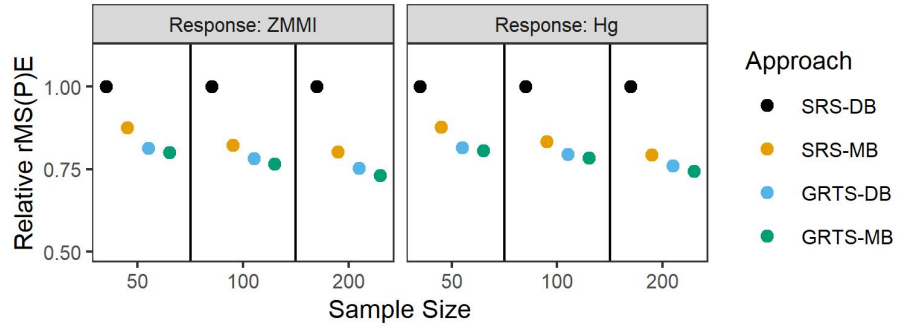


Figure 6: Relative rMS(P)E in the data study for the four sampling-inference combinations. The rows indicate the proportion of dependent error and the columns indicate the response type.

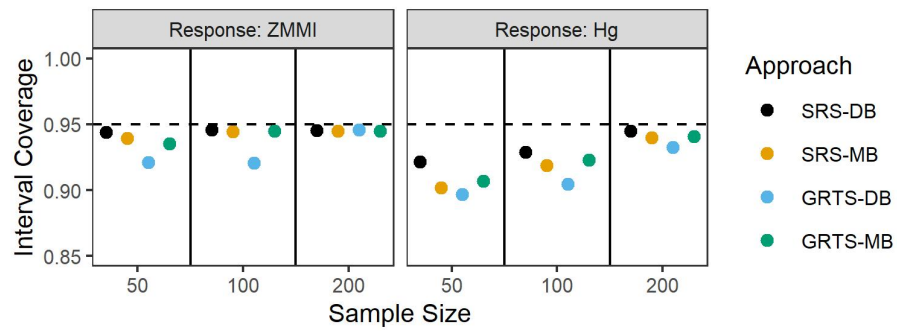


Figure 7: Interval coverage in the data study for the four sampling-inference combinations. The rows indicate the proportion of dependent error and the columns indicate the response type. The solid, black line represents 95% coverage.

data – GRTS-Design and GRTS-Model generally had similar rMS(P)E as their IRS counterparts when there was no spatial covariance and lower rMS(P)E than their IRS counterparts when there was spatial covariance. Second, the sampling decision (IRS vs GRTS) is most important when using a design-based analysis. Though GRTS-Model still had lower rMS(P)E than IRS-Model, the model-based analysis mitigated most of the rMS(P)E inefficiencies that result from the IRS samples lacking spatial balance. Third, as the strength of spatial covariance increases, the gap in rMS(P)E and MStdE between IRS-Design and the other sampling-inference combinations also increases, likely because IRS-Design is the only combination that ignores spatial locations in sampling and analysis. Fourth and finally, when the response was normal, interval coverage for all sampling-inference combinations was usually close to 95% for all sample sizes; when the response was lognormal, interval coverage for all sampling-inference combinations was usually between 90% and 95% and closest to 95% when $n = 200$.

AT LEAST HAVE DISCUSSION ABOUT MODEL BASED ASSUMPTIONS
AND MOVE VALIDITY COMMENTS TO RESULTS SECTION.

There are several benefits and drawbacks of the design-based and model-based approaches for finite population spatial data. Some we have discussed, but others we have not, and they are worthy of consideration in future research. Design-based approaches are often computationally efficient, while model-based approaches can be computationally burdensome, especially for likelihood-based estimation methods like REML that rely on inverting a covariance matrix. The design-based approach easily handles binary data through a straightforward application of the Horvitz-Thompson estimator. In contrast, analyzing binary data using a model-based approach generally requires a logistic mixed regression model, which can be challenging to estimate and interpret (Bolker et al., 2009). The design-based approach yields valid results because the sampling plan and

471 inclusion probabilities are specified directly by the researcher, while the model-
472 based approach may not yield valid results if the assumptions made do not
473 not accurately capture reality. The model-based approach, however, can more
474 naturally quantify the relationship between covariates (predictor variables) and
475 the response variable. The model-based approach also yields estimated spatial
476 covariance parameters, which help better understand the dependence structure
477 in the process in study. Model selection is also possible using model-based
478 approaches and criteria such as cross validation, likelihood ratio tests, or AIC
479 (Akaike, 1974). Model-based approaches are capable of more efficient small-area
480 estimation than design-based approaches by leveraging distributional assumptions
481 in areas with few observed units. Model-based approaches can also compute
482 unit-by-unit predictions at unobserved locations and use them to construct
483 informative visualizations like smoothed maps. Brus and De Gruijter (1997)
484 provide a more thorough discussion regarding the benefits and drawbacks of the
485 two approaches. In short, when deciding whether the design-based or model-
486 based approach is more appropriate to implement, the benefits and drawbacks of
487 each approach should be considered alongside the particular goals of the study.

488 **Acknowledgments**

489 We would like to thank the editors and anonymous reviewers for their
490 thoughtful comments which greatly improved the manuscript.

491 The views expressed in this manuscript are those of the authors and do not
492 necessarily represent the views or policies of the U.S. Environmental Protection
493 Agency or the National Oceanic and Atmospheric Administration. Any mention
494 of trade names, products, or services does not imply an endorsement by the
495 U.S. government, the U.S. Environmental Protection Agency, or the National
496 Oceanic and Atmospheric Administration. The U.S. Environmental Protection

497 Agency and National Oceanic and Atmospheric Administration do not endorse
498 any commercial products, services, or enterprises.

499 **Conflict of Interest Statement**

500 There are no conflicts of interest for any of the authors.

501 **Author Contribution Statement**

502 All authors conceived the ideas; All authors designed the methodology; MD
503 and MH performed the simulations and analyzed the data; MD and MH led the
504 writing of the manuscript; All authors contributed critically to the drafts and
505 gave final approval for publication.

506 **Data and Code Availability**

507 This manuscript has a supplementary **R** package that contains all of the
508 data and code used in its creation. The supplementary **R** package is hosted on
509 GitHub. Instructions for download at available at

510 <https://github.com/michaeldumelle/DvMsp>.

511 If the manuscript is accepted, this repository will be archived in Zenodo.

512 **Supporting Information**

513 In the supporting information, we provide tables of summary statistics for
514 all 36 simulation scenarios.

515 **References**

516 Akaike, H., 1974. A new look at the statistical model identification. IEEE
517 Transactions on Automatic Control 19, 716–723.

- 518 Barabesi, L., Franceschi, S., 2011. Sampling properties of spatial total
519 estimators under tessellation stratified designs. *Environmetrics* 22, 271–278.
- 520 Benedetti, R., Piersimoni, F., 2017. A spatially balanced design with proba-
521 bility function proportional to the within sample distance. *Biometrical Journal*
522 59, 1067–1084.
- 523 Benedetti, R., Piersimoni, F., Postiglione, P., 2017. Spatially balanced
524 sampling: A review and a reappraisal. *International Statistical Review* 85,
525 439–454.
- 526 Bolker, B.M., Brooks, M.E., Clark, C.J., Geange, S.W., Poulsen, J.R.,
527 Stevens, M.H.H., White, J.-S.S., 2009. Generalized linear mixed models: A
528 practical guide for ecology and evolution. *Trends in ecology & evolution* 24,
529 127–135.
- 530 Breiman, L., 2001. Random forests. *Machine Learning* 45, 5–32.
- 531 Brus, D., De Gruijter, J., 1997. Random sampling or geostatistical modelling?
532 Choosing between design-based and model-based sampling strategies for soil
533 (with discussion). *Geoderma* 80, 1–44.
- 534 Brus, D.J., 2021. Statistical approaches for spatial sample survey: Persistent
535 misconceptions and new developments. *European Journal of Soil Science* 72,
536 686–703.
- 537 Brus, D.J., DeGruijter, J.J., 1993. Design-based versus model-based esti-
538 mates of spatial means: Theory and application in environmental soil science.
539 *Environmetrics* 4, 123–152.
- 540 Chan-Golston, A.M., Banerjee, S., Handcock, M.S., 2020. Bayesian inference
541 for finite populations under spatial process settings. *Environmetrics* 31, e2606.
- 542 Chiles, J.-P., Delfiner, P., 1999. *Geostatistics: Modeling Spatial Uncertainty*.
543 John Wiley & Sons, New York.
- 544 Cicchitelli, G., Montanari, G.E., 2012. Model-assisted estimation of a spatial

- 545 population mean. *International Statistical Review* 80, 111–126.
- 546 Cooper, C., 2006. Sampling and variance estimation on continuous domains.
547 *Environmetrics* 17, 539–553.
- 548 Cressie, N., 1993. *Statistics for spatial data*. John Wiley & Sons.
- 549 De Gruijter, J., Ter Braak, C., 1990. Model-free estimation from spatial
550 samples: A reappraisal of classical sampling theory. *Mathematical Geology* 22,
551 407–415.
- 552 Diggle, P.J., Menezes, R., Su, T.-l., 2010. Geostatistical inference under
553 preferential sampling. *Journal of the Royal Statistical Society: Series C (Applied*
554 *Statistics)* 59, 191–232.
- 555 Dumelle, M., Kincaid, T.M., Olsen, A.R., Weber, M.H., 2022. *Spsurvey*:
556 *Spatial sampling design and analysis*.
- 557 Fix, E., Hodges, J.L., 1989. Discriminatory analysis. Nonparametric dis-
558 crimination: Consistency properties. *International Statistical Review/Revue*
559 *Internationale de Statistique* 57, 238–247.
- 560 Grafström, A., 2012. Spatially correlated poisson sampling. *Journal of*
561 *Statistical Planning and Inference* 142, 139–147.
- 562 Grafström, A., Lundström, N.L., 2013. Why well spread probability samples
563 are balanced. *Open Journal of Statistics* 3, 36–41.
- 564 Grafström, A., Lundström, N.L., Schelin, L., 2012. Spatially balanced
565 sampling through the pivotal method. *Biometrics* 68, 514–520.
- 566 Grafström, A., Matei, A., 2018. Spatially balanced sampling of continuous
567 populations. *Scandinavian Journal of Statistics* 45, 792–805.
- 568 Hansen, M.H., Madow, W.G., Tepping, B.J., 1983. An evaluation of model-
569 dependent and probability-sampling inferences in sample surveys. *Journal of the*
570 *American Statistical Association* 78, 776–793.
- 571 Harville, D.A., 1977. Maximum likelihood approaches to variance compo-

- 572 nent estimation and to related problems. *Journal of the American Statistical*
573 *Association* 72, 320–338.
- 574 Hofman, S.C., Brus, D., 2021. How many sampling points are needed to
575 estimate the mean nitrate-n content of agricultural fields? A geostatistical
576 simulation approach with uncertain variograms. *Geoderma* 385, 114816.
- 577 Horvitz, D.G., Thompson, D.J., 1952. A generalization of sampling with-
578 out replacement from a finite universe. *Journal of the American Statistical*
579 *Association* 47, 663–685.
- 580 Lohr, S.L., 2009. *Sampling: Design and analysis*. Nelson Education.
- 581 Patterson, H.D., Thompson, R., 1971. Recovery of inter-block information
582 when block sizes are unequal. *Biometrika* 58, 545–554.
- 583 Robertson, B., Brown, J., McDonald, T., Jaksons, P., 2013. BAS: Balanced
584 acceptance sampling of natural resources. *Biometrics* 69, 776–784.
- 585 Robertson, B., McDonald, T., Price, C., Brown, J., 2018. Halton iterative
586 partitioning: Spatially balanced sampling via partitioning. *Environmental and*
587 *Ecological Statistics* 25, 305–323.
- 588 Särndal, C.-E., Swensson, B., Wretman, J., 2003. *Model assisted survey*
589 *sampling*. Springer Science & Business Media.
- 590 Schabenberger, O., Gotway, C.A., 2017. *Statistical methods for spatial data*
591 *analysis*. CRC press.
- 592 Sen, A.R., 1953. On the estimate of the variance in sampling with varying
593 probabilities. *Journal of the Indian Society of Agricultural Statistics* 5, 127.
- 594 Sterba, S.K., 2009. Alternative model-based and design-based frameworks
595 for inference from samples to populations: From polarization to integration.
596 *Multivariate Behavioral Research* 44, 711–740.
- 597 Stevens, D.L., Olsen, A.R., 2003. Variance estimation for spatially balanced
598 samples of environmental resources. *Environmetrics* 14, 593–610.

- 599 Stevens, D.L., Olsen, A.R., 2004. Spatially balanced sampling of natural
600 resources. *Journal of the American Statistical Association* 99, 262–278.
- 601 USEPA, 2012. National lakes assessment 2012. [https://www.epa.gov/national-](https://www.epa.gov/national-aquatic-resource-surveys/national-results-and-regional-highlights-national-lakes-assessment)
602 aquatic-resource-surveys/national-results-and-regional-highlights-national-lakes-
603 assessment.
- 604 Ver Hoef, J., 2002. Sampling and geostatistics for spatial data. *Ecoscience* 9,
605 152–161.
- 606 Ver Hoef, J.M., 2008. Spatial methods for plot-based sampling of wildlife
607 populations. *Environmental and Ecological Statistics* 15, 3–13.
- 608 Ver Hoef, J.M., Temesgen, H., 2013. A comparison of the spatial linear
609 model to nearest neighbor (k-nn) methods for forestry applications. *PLOS ONE*
610 8, e59129.
- 611 Walvoort, D.J., Brus, D., De Gruijter, J., 2010. An r package for spatial
612 coverage sampling and random sampling from compact geographical strata by
613 k-means. *Computers & geosciences* 36, 1261–1267.
- 614 Wang, J.-F., Jiang, C.-S., Hu, M.-G., Cao, Z.-D., Guo, Y.-S., Li, L.-F., Liu, T.-
615 J., Meng, B., 2013. Design-based spatial sampling: Theory and implementation.
616 *Environmental Modelling & Software* 40, 280–288.
- 617 Wolfinger, R., Tobias, R., Sall, J., 1994. Computing gaussian likelihoods and
618 their derivatives for general linear mixed models. *SIAM Journal on Scientific*
619 *Computing* 15, 1294–1310.