

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 lie on fundamentally different foundations. In the design-based approach, inference depends on random sampling. In the model-based approach, inference depends 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 simulations and an analysis of real mercury concentration data. In the simulations, a variety of sample sizes, location layouts, dependence structures, and response types are considered. In the simulations and real data analysis, 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 simulations and mercury concentration data, we found that regardless of the strength of spatial dependence in the data, sampling plans that incorporate spatial locations (spatially balanced samples) generally outperform sampling plans that ignore spatial locations (non-spatially balanced samples). We also found that model-based approaches tend to

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

Preprint submitted to *Methods in Ecology and Evolution*

December 21, 2021

outperform design-based approaches, even when the data are skewed (and by consequence, the model-based distributional assumptions violated). The performance gap between the analysis approaches is small when spatially balanced samples are used but large when non-spatially balanced samples are used. This suggests that the sampling choice (whether to select a sample that is spatially balanced) 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; Model-based inference; Spatially balanced sampling; Spatial covariance;

1. Introduction

There are two general approaches for using data to make frequentist statistical inferences about a population: design-based and model-based. 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 of population units is called a sample. In the design-based approach, inferences about the underlying population are informed via a probabilistic process that randomly assigns some population units to be in the sample. Alternatively, in the model-based approach, inferences

are made from specific assumptions about the underlying process generating the data. Each paradigm has a deep historical context (Sterba, 2009) and its own set of benefits and drawbacks (Hansen et al., 1983).

Though the design-based and model-based approaches apply to statistical inference in a broad sense, we focus on comparing these approaches for spatial data. We define spatial data as data that incorporates the specific locations of the population units into either the design or estimation process. De Gruijter and Ter Braak (1990) give an early comparison of design-based and model-based approaches for spatial data, quashing the belief that design-based approaches could not be used for spatially correlated data. Since then, there have been several general comparisons between design-based and model-based approaches for spatial data (Brus and De Gruijter, 1997; Brus, 2021; Ver Hoef, 2002, 2008; Wang et al., 2012). Cooper (2006) reviews the two approaches in an ecological context before introducing a “model-assisted” variance estimator that combines aspects from each approach. In addition to Cooper (2006), there has been substantial research and development into estimators that use both design and model-based principles (see e.g., Sterba (2009) and Cicchitelli and Montanari (2012), and see Chan-Golston et al. (2020) for a Bayesian approach).

Certainly comparisons between design-based and model-based approaches to spatial data have been studied. But no numerical comparison has been made between design-based approaches that incorporate spatial information and model-based approaches. In this manuscript, we compare design-based approaches that incorporate spatial information to model-based approaches for finite population spatial data. A finite population contains a finite number of population units (we assume the finite number is known); an example is lakes (treated as a whole with the lake centroid representing location) in the contiguous United States. Though we manuscript focuses on finite populations, these comparisons generalize to

82 infinite populations as well. An infinite population contains an infinite number of
83 population units; an example is locations within a single lake. In this manuscript
84 we assume the number of finite population units is known

85 The rest of the manuscript is organized as follows. In Section 1.1, we
86 introduce and provide relevant background for the design-based and model-based
87 approaches to finite population spatial data. In Section 2, we describe how
88 we compare performance of the approaches with a simulation study and an
89 analysis of real data that contains mercury concentration in lakes located in the
90 contiguous United States. In Section 3, we present results from the simulation
91 study and the mercury concentration analysis. And in Section 4, we end with a
92 discussion and provide directions for future research.

93 *1.1. Background*

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

98 *1.1.1. Comparing Design-Based and Model-Based Approaches*

99 The design-based approach assumes the population is fixed. Randomness
100 is incorporated via the selection of units in a sampling frame. A sampling
101 frame is the set of all units available to be sampled. Units from the sampling
102 frame are selected as part of the sample according to a sampling design, which
103 assigns a positive probability of inclusion (inclusion probability) to each unit
104 from the sampling frame. These inclusion probabilities are later used to analyze
105 data. Some examples of commonly used sampling designs include simple random
106 sampling, stratified random sampling, and cluster sampling.

107 When sampling designs incorporate spatial locations into sampling, we call
108 the resulting samples “spatially balanced.” One approach to selecting spatially

109 balanced samples is the Generalized Random Tessellation Stratified (GRTS)
110 algorithm (Stevens and Olsen, 2004), which we discuss in more detail in Section
111 1.1.2. When sampling designs do not incorporate spatial locations into sampling,
112 we call the resulting samples “non-spatially balanced.”

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

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

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

149 1.1.2. *Spatially Balanced Design and Analysis*

150 We previously mentioned that the design-based approach can be used to
 151 select spatially balanced samples (samples that incorporate spatial locations of
 152 the population units). Spatially balanced samples are useful because paramete-
 153 ter estimates from these samples tend to vary less than parameter estimates
 154 from samples that are not spatially balanced (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
 160 on Voronoi polygons (Dirichlet Tessellations). After the GRTS algorithm was
 161 developed, several other spatially balanced sampling algorithms emerged, such as
 162 the Local Pivotal Method (Grafström et al., 2012; Grafström and Matei, 2018),

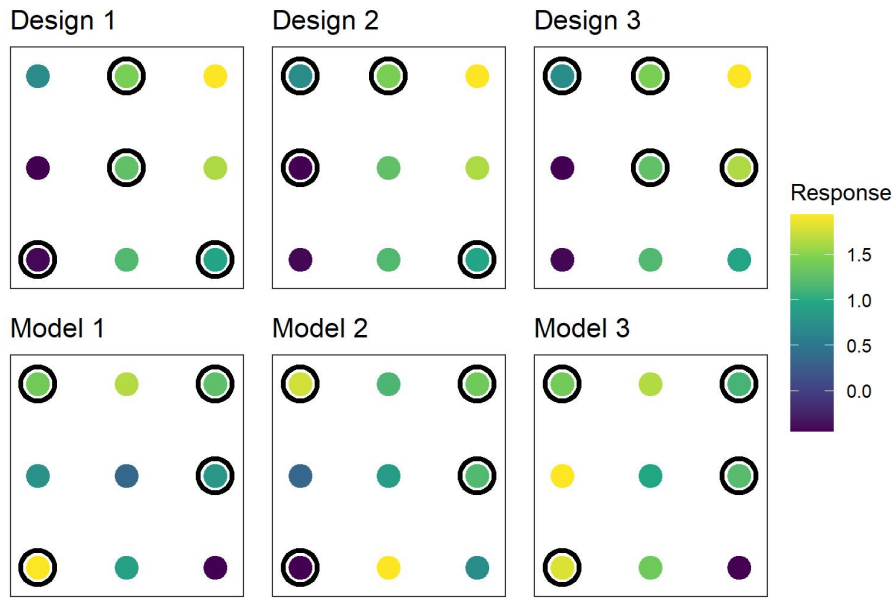


Figure 1: A visual comparison of the design-based and model-based approaches. In the top row, there is one fixed population with nine population units and three random samples of size four (points circled are those sampled). The response values at each site are fixed, but we obtain different estimates for the mean response in each random sample. In the bottom row, there are three realizations of the same data-generating stochastic process that are all sampled at the same four locations. The data-generating stochastic process has a single mean, but the mean of the nine population units is different in each of the three realizations

163 Spatially Correlated Poisson Sampling (Grafström, 2012), Balanced Acceptance
 164 Sampling (Robertson et al., 2013), Within-Sample-Distance Sampling (Benedetti
 165 and Piersimoni, 2017), and Halton Iterative Partitioning Sampling (Robertson
 166 et al., 2018). In this manuscript, we select spatially balanced samples using the
 167 Generalized Random Tessellation Stratified (GRTS) algorithm because it has sev-
 168 eral attractive properties. More specifically, the GRTS algorithm accommodates
 169 finite and infinite sampling frames, equal, unequal, and proportional (to size) in-
 170 clusion probabilities, legacy (historical) sampling (Foster et al., 2017), a minimum
 171 distance between units in a sample, and replacement units (replacement units are
 172 population units that can be sampled when a population unit originally selected
 173 can no longer be sampled). The GRTS algorithm selects samples by utilizing a
 174 particular mapping between two-dimensional and one-dimensional space that
 175 preserves proximity relationships. Via this mapping, units in two-dimensional
 176 space are partitioned using a hierarchical address. This hierarchical address is
 177 used to map population units to a one-dimensional line. On the one dimensional
 178 line, each population unit’s line length equals its inclusion probability. Then, a
 179 systematic sample of population units is selected on the line and mapped back
 180 to two-dimensional space, yielding the desired sample. Stevens and Olsen (2004)
 181 provide more technical details.

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

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

182 where Z_i is the value of the i th population unit in the sample and π_i is the
 183 inclusion probability of the i th population unit in the sample. An estimate of

the population mean is obtained by dividing $\hat{\tau}_{ht}$ by N , the number of population units.

It is also important to quantify uncertainty $\hat{\tau}_{ht}$. Horvitz and Thompson (1952) and Sen (1953) provide variance estimators for $\hat{\tau}_{ht}$, but these estimators have two drawbacks. First, they rely on calculating π_{ij} , the probability that population unit i and population unit j are both in the sample – this quantity can be challenging if not impossible to calculate analytically. Second, these estimators ignore the spatial locations of the population units. To address these two drawbacks simultaneously, Stevens and Olsen (2003) proposed the local neighborhood variance estimator. The local neighborhood variance estimator does not rely on π_{ij} and incorporates spatial locations – for technical details see Stevens and Olsen (2003). Stevens and Olsen (2003) show the local neighborhood variance estimator tends to reduce the estimated variance of $\hat{\tau}$ and yield narrower confidence intervals compared to variance estimators that ignore spatial locations.

1.1.3. Finite Population Block Kriging

Finite Population Block Kriging (FPBK) is a model-based approach that expands the geostatistical Kriging framework to the finite population setting (Ver Hoef, 2008). Instead of developing inference based on a specific sampling design, we assume the data are generated by a spatial stochastic process. We summarize some of the basic principles of FBPK next (for more technical details, see Ver Hoef (2008)). Let $\mathbf{z} \equiv \{z(s_1), z(s_2), \dots, z(s_N)\}$ be an $N \times 1$ response vector at locations s_1, s_2, \dots, s_N that can be measured at the N population units. Suppose we want to use a sample to predict some linear function of the response 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 one). 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

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.

FBPK 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 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 (for a thorough list of spatial covariance functions, see Cressie (1993). 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 quantifying the variability that is dependent (coarse-scale), σ_2^2 is the variance parameter quantifying the variability that is independent (fine-scale), ϕ is the range parameter measuring the distance-decay rate of the covariance, and $h_{i,j}$ is the Euclidean distance between population

units i and j . The proportion of variability attributable to dependent random error is $\sigma_1^2/(\sigma_1^2 + \sigma_2^2)$. Similarly, the proportion of variability attributable to independent random error is $\sigma_2^2/(\sigma_1^2 + \sigma_2^2)$. Finally we note that σ_1^2 and σ_2^2 are often called the partial sill and nugget, respectively.

With the above model formulation, the Best Linear Unbiased Predictor (BLUP) for $f(\mathbf{b}'\mathbf{z})$ and its prediction variance can be computed. While details of the derivation are in Ver Hoef (2008), we note here that the predictor and its variance are both moment-based, meaning that they do not rely on any distributional assumptions.

Other approaches, such as k-nearest-neighbors (Fix and Hodges, 1989; Ver Hoef and Temesgen, 2013) and random forests (Breiman, 2001), among others, could also be used to obtain predictions for a mean or total from spatially correlated responses of a finite population. Compared to the k-nearest-neighbors and random forest approach, we prefer FBPK because it is model-based and relies on theoretically-based variance estimators leveraging the model’s spatial covariance structure, whereas k-nearest-neighbors and random forests use ad-hoc variance estimators (Ver Hoef and Temesgen, 2013). Additionally, Ver Hoef and Temesgen (2013) studied compared FBPK, k-nearest-neighbors, and random forests in a variety of spatial data contexts, and FBPK tended to perform best.

2. Materials and Methods

2.1. Simulation Study

We used a simulation study to investigate performance of four sampling-analysis combinations: IRS sampling with a design-based analysis, called “IRS-Design”; IRS sampling with a model-based analysis, called “IRS-Model”; GRTS sampling with a design-based analysis, called “GRTS-Design”; GRTS sampling with a model-based analysis, called “GRTS-Model”. These combinations are also

provided in Table 1.

| | Design | Model |
|------|-------------|------------|
| IRS | IRS-Design | IRS-Model |
| GRTS | GRTS-Design | GRTS-Model |

Table 1: Sampling-analysis combinations in the simulation study. The rows give the two types of sampling designs and the columns give the two types of analyses.

Performance for the four sampling-analysis combinations was evaluated in 36 different simulation scenarios. The 36 scenarios resulted from the crossing of three sample sizes, two location layouts, two response types, and three proportions of dependent random error. 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 (of the population units) 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 lognormal. For the normal response type, the response was simulated using mean-zero random errors with the exponential covariance (Equation 3) for varying proportions of dependent random error. The proportion of dependent random error is represented by $\sigma_1^2/(\sigma_1^2 + \sigma_2^2)$, where σ_1^2 and σ_2^2 are the dependent random error variance (partial sill) and independent random error variance (nugget), respectively, from Equation 3. The total variance, $\sigma_1^2 + \sigma_2^2$, was always 2. The range was always $\sqrt{2}/3$, which means that the correlation in the dependent random error decayed to nearly zero at the largest possible distance between two population units in the domain. For the lognormal 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 lognormal random variable whose total variance is 2. The lognormal responses were used to evaluate performance of the sampling-

analysis approaches for data that were skewed (i.e., not normal).

| | | | |
|-------------------------------|--------|-----------|-----|
| Sample Size (n) | 50 | 100 | 200 |
| Location Layout | Random | Gridded | - |
| Proportion of Dependent Error | 0 | 0.5 | 0.9 |
| Response Type | Normal | Lognormal | - |

Table 2: Simulation scenario options. All combinations of sample size, location layout, response type, and proportion of dependent random error composed the 36 simulation scenarios. In each simulation scenario, the total variance was 2.

In each of the 36 simulation scenarios, there were 2000 independent simulation trials. In each 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 confidence (design-based) or prediction (model-based) intervals. Then we recorded the bias, squared error, and interval coverage for all sampling-analysis combinations. After all 2000 trials, we summarized the long-run performance of the combinations by calculating average bias, RMS(P)E (root-mean-squared error for the design-based approaches and root-mean-squared-prediction error for the model-based approaches), and the proportion of times the true mean is contained in its 95% interval. The GRTS algorithm and the local neighborhood variance estimator are available in the **R** package `spsurvey` (Dumelle et al., 2021). FPBK is available in the `sptotal` **R** package (Higham et al., 2021) and covariance parameters were estimated using Restricted Maximum Likelihood (Harville, 1977; Patterson and Thompson, 1971; Wolfinger et al., 1994).

2.2. Application

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. We will use data from the 2012 National Lakes Assessment (NLA), which measures various aspects of lake health and water quality (USEPA, 2012). Specifically,

we will analyze mercury concentration in lakes. Although we can calculate the true mean mercury concentration values for the 986 lakes from the 2012 NLA, we will explore whether or not we obtain an adequately precise estimate for the realized mean mercury concentration if we sample only 100 of the 986 lakes. For each of the four familiar sampling-analysis 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% confidence (design-based) or prediction (model-based) intervals from this sample of 100 lakes, which we compare to the actual mean from all 986 lakes.

3. Results

3.1. Simulation Study

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

Fig. 2 shows the relative rMS(P)E of the four approaches from Table 1 using the random location layout with “IRS-Design” as the baseline. The relative rMS(P)E is defined as

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

When there is no spatial correlation (Fig. 2, “Prop DE: 0” row), the four sampling-analysis combinations have approximately equal rMS(P)E. So using the GRTS sampling plan or a model-based analysis does not result in much, if any, loss in efficiency compared to IRS-Design when there is no spatial correlation. When there is spatial correlation (Fig. 2, “Prop DE: 0.5” and “Prop DE: 0.9” rows), GRTS-Model tends to perform best, followed by GRTS-Design, IRS-Model, and finally IRS-Design, though the difference in relative rMS(P)E among

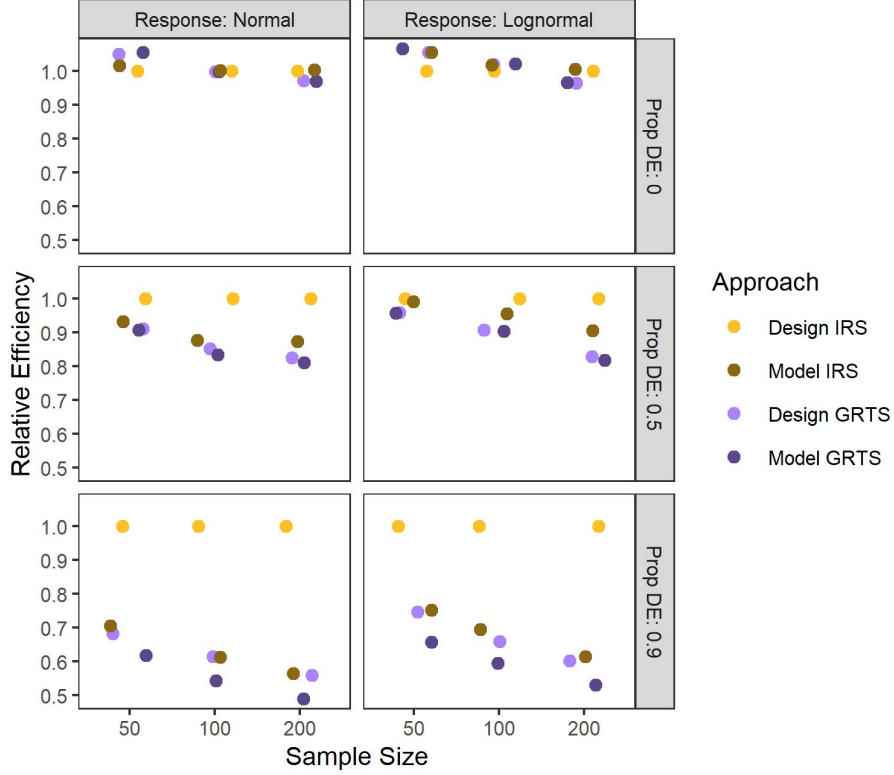


Figure 2: Relative rMS(P)E in the simulation study for the four sampling-analysis combinations. The rows indicate the proportion of dependent error and the columns indicate the response type.

312 GRTS-Model, GRTS-Design, and IRS-Model is relatively small. As the strength
 313 of spatial correlation increases, the gap in rMS(P)E between IRS-Design and the
 314 other sampling-analysis combinations widens. Finally we note that when there
 315 is spatial correlation, IRS-Model outperforms IRS-Design by a large margin,
 316 suggesting that the poor design properties of IRS are largely mitigated by the
 317 model-based analysis. These conclusions are similar to those observed in the grid
 318 location layout, so we omit a grid location layout figure here. Tables for rMS(P)E
 319 in all 36 simulation scenarios are provided in the supporting information.

320 We also studied 95% interval coverage among the sampling-analysis com-
 321 binations. The design-based confidence intervals and model-based prediction

intervals were constructed using the normal distribution. Justification for this comes from the asymptotic normality of means via the Central Limit Theorem. Fig. 3 shows the 95% interval coverage for each of the four sampling-analysis combinations in the random location layout. Within each scenario, the sampling-analysis combinations tend to have fairly similar interval coverage. Coverage in the normal response scenarios was usually near 95%, while coverage in the lognormal response scenarios varied from 90% to 95% but increased with the sample size. At a sample size of 200, all four sampling-analysis combinations had approximately 95% interval coverage in both response scenarios for all dependent error proportions. These 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.

3.2. Application

Fig. 4 shows a map and histogram of mercury concentration. The map shows mercury concentration exhibits some spatial patterning, with high mercury concentrations in lakes 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. 4 also shows mercury concentration's empirical semivariogram. The empirical semivariogram can be used as a tool to visualize spatial dependence. It quantifies the halved squared differences (semivariance) among mercury concentration at different distances apart. When a process is spatially correlated (has spatial dependence), the semivariance tends to be smaller at small distances and larger at large distances. The empirical semivariogram in Fig. 4 suggests that mercury concentration exhibits spatial dependence. Lastly we note that the realized mean mercury concentration in the 986 NLA lakes is

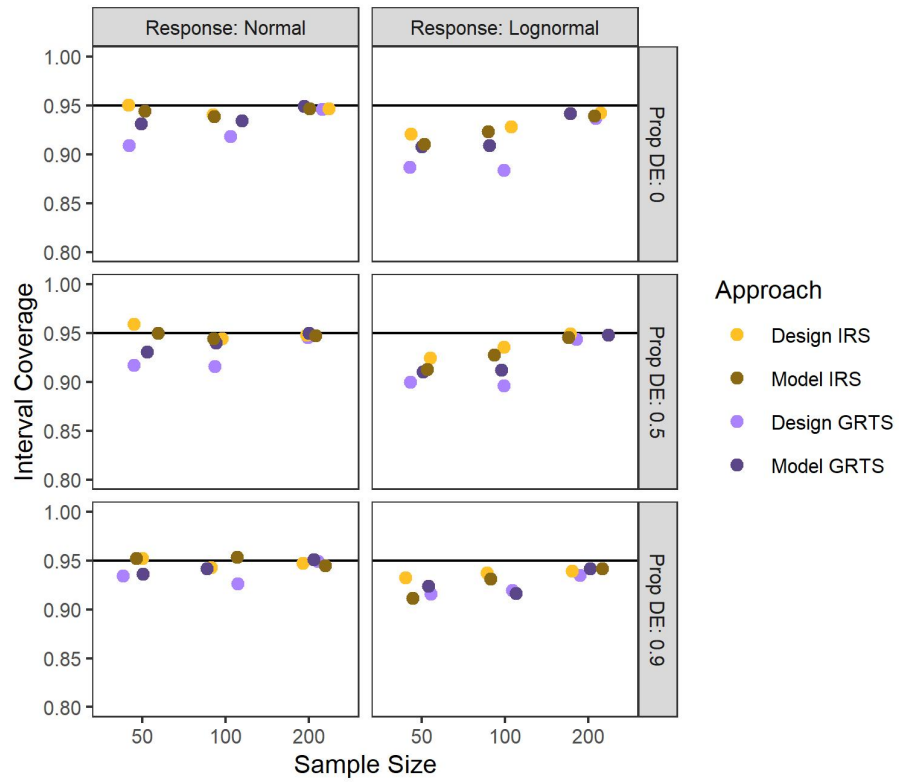


Figure 3: Interval coverage in the simulation study for the four sampling-analysis combinations. The rows indicate the proportion of dependent error and the columns indicate the response type. The solid, black line in each plot represents 95% coverage.

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 its standard error using design-based and model-based approaches. For the model-based analyses, the exponential covariance was used. Table 3 shows the results from these analyses. For all four sampling-analysis combinations, the true realized mean mercury concentration is within the bounds of the 95% confidence (design-based) or prediction (model-based) intervals. Though we should not generalize these results to other samples from these data, we do note a couple of patterns. The design-based IRS analysis shows the largest standard error: a likely reason is that this is the only approach that does not incorporate any spatial information regarding mercury concentration. Both analyses using GRTS sampling have lower standard errors than both analyses using IRS sampling. We expect that these patterns are consistent with other samples from these data because mercury concentration exhibits spatial patterning, so a spatially balanced sample should usually yield a lower standard error.

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

Table 3: For each sampling-analysis combination (Approach), 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 the sample of 100 lakes in the NLA data. The true mean concentration of all 986 lakes in the NLA data is 103.2 ng / g.

4. Discussion

The design-based and model-based approaches to statistical inference are fundamentally different paradigms that can be used to analyze data. The design-based approach incorporates randomness through sampling to estimate

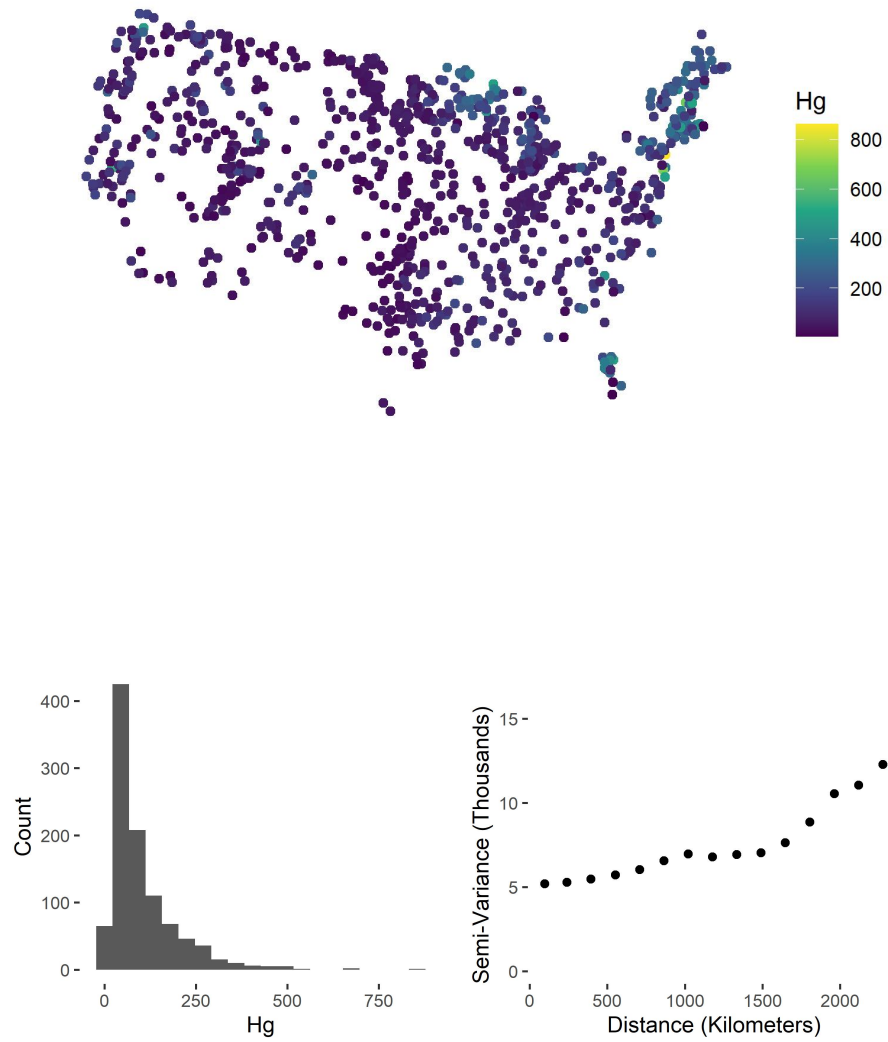


Figure 4: Mercury concentration visualizations for the population (Hg) for 986 lakes in the NLA data. A spatial layout is in the top row, a histogram is in the bottom row and left column, and an empirical semivariogram is in the bottom row and right column.

369 population parameters. The model-based approach incorporates randomness
 370 through distributional assumptions to predict realized values of a random process.
 371 Though these approaches have often been compared in the literature both from
 372 theoretical and analytical perspectives, our contribution lies in studying them
 373 in a spatial context while implementing spatially balanced sampling. Aside
 374 from the theoretical differences described, a few analytical findings from the
 375 simulation study are particularly notable. First, the sampling decision (GRTS
 376 vs IRS) is most important when using a design-based analysis. Though GRTS-
 377 Model still outperformed IRS-Model, the model-based analysis mitigated much
 378 of the inefficiency of the IRS sample. Second, independent of the analysis
 379 approach, we found no reason to prefer IRS over GRTS for sampling spatial data
 380 – GRTS-Design and GRTS-Model generally performed at least as well as their IRS
 381 counterparts when there was no spatial correlation and noticeably better than
 382 their IRS counterparts when there was spatial correlation. Third, as the strength
 383 of spatial correlation increases, the gap in rMS(P)E between IRS-Design and the
 384 other sampling-analysis combinations also increases. Fourth and finally, when
 385 the response was normal, interval coverage for all sampling-analysis combinations
 386 was very close to 95% for all sample sizes; when the response was lognormal,
 387 interval coverage for all sampling and analysis was between 90% and 95% and
 388 closest to 95% when $n = 200$.

389 There are several benefits and drawbacks of the design-based and model-
 390 based approaches for finite population spatial data. Some we have discussed,
 391 but others we have not, and they are worthy of consideration in future research.
 392 Design-based approaches are often computationally efficient, while model-based
 393 approaches can be computationally burdensome, especially for likelihood-based
 394 estimation methods like REML that rely on inverting a covariance matrix. The
 395 design-based approach also more naturally handles binary data, free from the

396 more complicated logistic regression framework commonly used to analyze binary
397 data in a model-based approach. The model-based approach, however, can more
398 naturally quantify the relationship between covariates (predictor variables) and
399 response variable. The model-based approach also yields estimated spatial
400 covariance parameters, which help better understand the dependence structure
401 in the process of study. Model selection is also possible using model-based
402 approaches and criteria such as cross validation, likelihood ratio tests, or AIC
403 (Akaike, 1974). Model-based approaches are capable of more efficient small-area
404 estimation than design-based approaches by leveraging distributional assumptions
405 in areas with few observed sites. Model-based approaches can also compute site-
406 by-site predictions at unobserved locations and use them to construct informative
407 visualizations. The benefits and drawbacks of the design-based and model-based
408 approaches should be considered alongside the particular goals of a study when
409 deciding which approach is most appropriate to implement.

410 **Acknowledgments**

411 The views expressed in this manuscript are those of the authors and do not
412 necessarily represent the views or policies of the U.S. Environmental Protection
413 Agency or the National Oceanic and Atmospheric Administration. Any mention
414 of trade names, products, or services does not imply an endorsement by the
415 U.S. government, the U.S. Environmental Protection Agency, or the National
416 Oceanic and Atmospheric Administration. The U.S. Environmental Protection
417 Agency and National Oceanic and Atmospheric Administration do not endorse
418 any commercial products, services, or enterprises.

419 **Conflict of Interest Statement**

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

421 **Data and Code Availability**

422 This manuscript has a supplementary R package that contains all of the
423 data and code used in its creation. The supplementary R package is hosted on
424 GitHub. Instructions for download are available at
425 <https://github.com/michaeldumelle/DvMsp>.

426 **Supporting Information**

427 In the supporting information, we provide tables presenting summary statis-
428 tics for all 36 simulation scenarios.

429 **Author Contributions**

430 All authors conceived the ideas; All authors designed methodology; MD and
431 MH performed the simulations and analyzed the data; MD and MH led the
432 writing of the manuscript; All authors contributed critically to the drafts and
433 gave final approval for publication.

434 **References**

- 435 Akaike, H., 1974. A new look at the statistical model identification. IEEE
436 Transactions on Automatic Control 19, 716–723.
- 437 Barabesi, L., Franceschi, S., 2011. Sampling properties of spatial total
438 estimators under tessellation stratified designs. Environmetrics 22, 271–278.
- 439 Benedetti, R., Piersimoni, F., 2017. A spatially balanced design with proba-
440 bility function proportional to the within sample distance. Biometrical Journal
441 59, 1067–1084.
- 442 Benedetti, R., Piersimoni, F., Postiglione, P., 2017. Spatially balanced
443 sampling: A review and a reappraisal. International Statistical Review 85,
444 439–454.

445 Breiman, L., 2001. Random forests. *Machine Learning* 45, 5–32.

446 Brus, D., De Gruijter, J., 1997. Random sampling or geostatistical modelling?
 447 Choosing between design-based and model-based sampling strategies for soil
 448 (with discussion). *Geoderma* 80, 1–44.

449 Brus, D.J., 2021. Statistical approaches for spatial sample survey: Persistent
 450 misconceptions and new developments. *European Journal of Soil Science* 72,
 451 686–703.

452 Chan-Golston, A.M., Banerjee, S., Handcock, M.S., 2020. Bayesian inference
 453 for finite populations under spatial process settings. *Environmetrics* 31, e2606.

454 Chiles, J.-P., Delfiner, P., 1999. *Geostatistics: Modeling Spatial Uncertainty*.
 455 John Wiley & Sons, New York.

456 Cicchitelli, G., Montanari, G.E., 2012. Model-assisted estimation of a spatial
 457 population mean. *International Statistical Review* 80, 111–126.

458 Cooper, C., 2006. Sampling and variance estimation on continuous domains.
 459 *Environmetrics* 17, 539–553.

460 Cressie, N., 1993. *Statistics for spatial data*. John Wiley & Sons.

461 De Gruijter, J., Ter Braak, C., 1990. Model-free estimation from spatial
 462 samples: A reappraisal of classical sampling theory. *Mathematical Geology* 22,
 463 407–415.

464 Diggle, P.J., Menezes, R., Su, T.-l., 2010. Geostatistical inference under
 465 preferential sampling. *Journal of the Royal Statistical Society: Series C (Applied*
 466 *Statistics)* 59, 191–232.

467 Dumelle, M., Kincaid, T.M., Olsen, A.R., Weber, M.H., 2021. *Spsurvey*:
 468 *Spatial sampling design and analysis*.

469 Fix, E., Hodges, J.L., 1989. Discriminatory analysis. Nonparametric dis-
 470 crimination: Consistency properties. *International Statistical Review/Revue*
 471 *Internationale de Statistique* 57, 238–247.

472 Foster, S.D., Hosack, G.R., Lawrence, E., Przeslawski, R., Hedge, P., Caley,
473 M.J., Barrett, N.S., Williams, A., Li, J., Lynch, T., others, 2017. Spatially
474 balanced designs that incorporate legacy sites. *Methods in Ecology and Evolution*
475 8, 1433–1442.

476 Grafström, A., 2012. Spatially correlated poisson sampling. *Journal of*
477 *Statistical Planning and Inference* 142, 139–147.

478 Grafström, A., Lundström, N.L., 2013. Why well spread probability samples
479 are balanced. *Open Journal of Statistics* 3, 36–41.

480 Grafström, A., Lundström, N.L., Schelin, L., 2012. Spatially balanced
481 sampling through the pivotal method. *Biometrics* 68, 514–520.

482 Grafström, A., Matei, A., 2018. Spatially balanced sampling of continuous
483 populations. *Scandinavian Journal of Statistics* 45, 792–805.

484 Hansen, M.H., Madow, W.G., Tepping, B.J., 1983. An evaluation of model-
485 dependent and probability-sampling inferences in sample surveys. *Journal of the*
486 *American Statistical Association* 78, 776–793.

487 Harville, D.A., 1977. Maximum likelihood approaches to variance compo-
488 nent estimation and to related problems. *Journal of the American Statistical*
489 *Association* 72, 320–338.

490 Higham, M., Ver Hoef, J., Frank, B., Dumelle, M., 2021. Sptotal: Predicting
491 totals and weighted sums from spatial data.

492 Horvitz, D.G., Thompson, D.J., 1952. A generalization of sampling with-
493 out replacement from a finite universe. *Journal of the American Statistical*
494 *Association* 47, 663–685.

495 Lohr, S.L., 2009. *Sampling: Design and analysis*. Nelson Education.

496 Patterson, H.D., Thompson, R., 1971. Recovery of inter-block information
497 when block sizes are unequal. *Biometrika* 58, 545–554.

498 Robertson, B., Brown, J., McDonald, T., Jaksons, P., 2013. BAS: Balanced

499 acceptance sampling of natural resources. *Biometrics* 69, 776–784.

500 Robertson, B., McDonald, T., Price, C., Brown, J., 2018. Halton iterative
501 partitioning: Spatially balanced sampling via partitioning. *Environmental and*
502 *Ecological Statistics* 25, 305–323.

503 Särndal, C.-E., Swensson, B., Wretman, J., 2003. Model assisted survey
504 sampling. Springer Science & Business Media.

505 Schabenberger, O., Gotway, C.A., 2017. Statistical methods for spatial data
506 analysis. CRC press.

507 Sen, A.R., 1953. On the estimate of the variance in sampling with varying
508 probabilities. *Journal of the Indian Society of Agricultural Statistics* 5, 127.

509 Sterba, S.K., 2009. Alternative model-based and design-based frameworks
510 for inference from samples to populations: From polarization to integration.
511 *Multivariate Behavioral Research* 44, 711–740.

512 Stevens, D.L., Olsen, A.R., 2003. Variance estimation for spatially balanced
513 samples of environmental resources. *Environmetrics* 14, 593–610.

514 Stevens, D.L., Olsen, A.R., 2004. Spatially balanced sampling of natural
515 resources. *Journal of the American Statistical Association* 99, 262–278.

516 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)
517 [aquatic-resource-surveys/national-results-and-regional-highlights-national-lakes-](https://www.epa.gov/national-aquatic-resource-surveys/national-results-and-regional-highlights-national-lakes-assessment)
518 [assessment](https://www.epa.gov/national-aquatic-resource-surveys/national-results-and-regional-highlights-national-lakes-assessment).

519 Ver Hoef, J., 2002. Sampling and geostatistics for spatial data. *Ecoscience* 9,
520 152–161.

521 Ver Hoef, J.M., 2008. Spatial methods for plot-based sampling of wildlife
522 populations. *Environmental and Ecological Statistics* 15, 3–13.

523 Ver Hoef, J.M., Temesgen, H., 2013. A comparison of the spatial linear
524 model to nearest neighbor (k-nn) methods for forestry applications. *PLOS ONE*
525 8, e59129.

526 Wang, J.-F., Jiang, C.-S., Hu, M.-G., Cao, Z.-D., Guo, Y.-S., Li, L.-F., Liu, T.-
527 J., Meng, B., 2013. Design-based spatial sampling: Theory and implementation.
528 Environmental Modelling & Software 40, 280–288.

529 Wang, J.-F., Stein, A., Gao, B.-B., Ge, Y., 2012. A review of spatial sampling.
530 Spatial Statistics 2, 1–14.

531 Wolfinger, R., Tobias, R., Sall, J., 1994. Computing gaussian likelihoods and
532 their derivatives for general linear mixed models. SIAM Journal on Scientific
533 Computing 15, 1294–1310.