

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

^aUnited States Environmental Protection Agency, 200 SW 35th St, Corvallis, Oregon, 97333

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

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

^dOregon 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 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 ap-

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

Preprint submitted to *Methods in Ecology and Evolution*

April 6, 2022

proaches tend to outperform design-based approaches, even when the data are skewed (and by consequence, the model-based distributional assumptions violated). The performance gap between these 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; 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

54 some population units to be in the sample (e.g., random sampling). Alterna-
 55 tively, in the model-based approach, inference relies on distributional assump-
 56 tions about the underlying stochastic process that generated the sample. Each
 57 paradigm has a deep historical context (Sterba, 2009) and its own set of benefits
 58 and drawbacks (Hansen et al., 1983).

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

74 Certainly comparisons between design-based and model-based approaches
 75 have been studied in spatial contexts. But no numerical comparison has been
 76 made between design-based approaches that incorporate spatial locations into
 77 sampling and analysis and model-based approaches. In this manuscript, we
 78 compare design-based approaches that incorporate spatial locations into sam-
 79 pling and analysis to model-based approaches for finite population spatial data.
 80 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 here we focus on finite populations, the comparisons we discuss generalize to infinite populations as well. An infinite population contains an infinite number of population units; an example is locations within a single lake.

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

1.1. Background

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

1.1.1. Comparing Design-Based and Model-Based Approaches

The design-based approach assumes the population is fixed. Randomness is incorporated via the selection of population units according to a sampling design. A sampling design assigns a positive probability of inclusion (inclusion probability) in the sample to each population unit. These inclusion probabilities are later used to estimate population parameters. Some examples of commonly used sampling designs include simple random 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
 116 (e.g., a population mean). Treating the data as fixed and incorporating ran-
 117 domness through the sampling design yields estimators having very few other
 118 assumptions. Confidence intervals for these types of estimators are typically
 119 derived using limiting arguments that incorporate all possible samples. Sample
 120 means, for example, are asymptotically normal (Gaussian) by the Central Limit
 121 Theorem (under some assumptions). If we repeatedly select samples from the
 122 population, then 95% of all 95% confidence intervals constructed from a pro-
 123 cedure with appropriate coverage will contain the true fixed population mean.
 124 Särndal et al. (2003) and Lohr (2009) provide thorough reviews of the design-
 125 based approach.

126 The model-based approach assumes the sample is a random realization of
 127 a data-generating stochastic process. Randomness is formally incorporated
 128 through distributional assumptions on this process. Strictly speaking, random-
 129 ness need not be incorporated through random sampling, though Diggle et al.
 130 (2010) warn against preferential sampling. Preferential sampling occurs when
 131 the process generating the data locations and the process being modeled are
 132 not independent of one another. To guard against preferential sampling, model-
 133 based approaches often still implement some form of random sampling. When

134 model-based approaches implement random sampling, the inclusion probabil-
 135 ities are ignored when analyzing the sample (in contrast to the design-based
 136 approach, which relies on these inclusion probabilities to analyze the sample).

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

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 (samples that incorporate spatial locations of
 153 the population units). Spatially balanced samples are useful because param-
 154 eter estimates from these samples tend to vary less than parameter estimates
 155 from samples that are not spatially balanced (Barabesi and Franceschi, 2011;
 156 Benedetti et al., 2017; Grafström and Lundström, 2013; Robertson et al., 2013;
 157 Stevens and Olsen, 2004; Wang et al., 2013). The first spatially balanced sam-
 158 pling algorithm to see widespread use was the Generalized Random Tessellation
 159 Stratified (GRTS) algorithm (Stevens and Olsen, 2004). To quantify the spatial
 160 balance of a sample, Stevens and Olsen (2004) proposed loss metrics based on

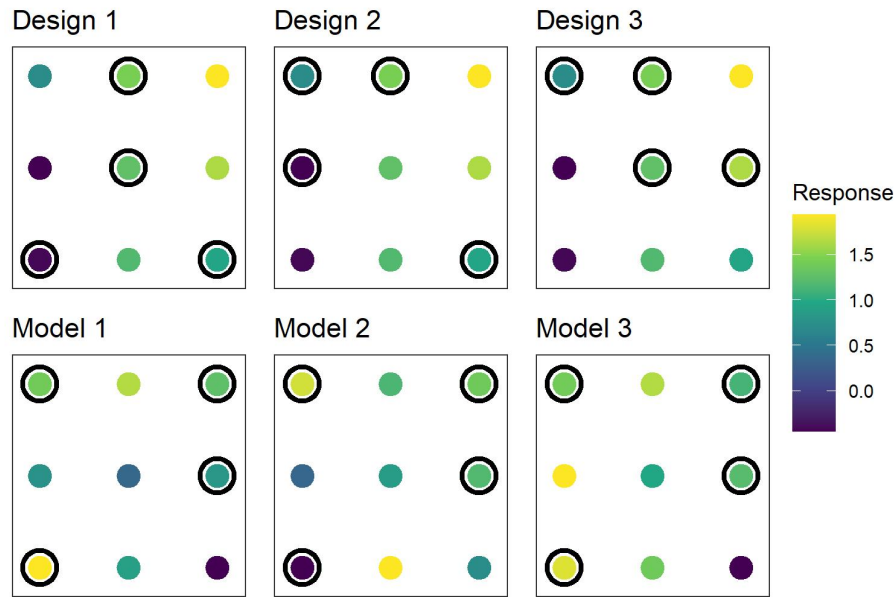


Figure 1: A visual comparison of the design-based and model-based approaches. In the top row, the design-based approach is highlighted. There is one fixed population with nine population units and three random samples of size four (points circled are those sampled). The response values at each site are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

161 Voronoi polygons (Dirichlet Tessellations). After the GRTS algorithm was de-
 162 veloped, several other spatially balanced sampling algorithms emerged, includ-
 163 ing the Local Pivotal Method (Grafström et al., 2012; Grafström and Matei,
 164 2018), Spatially Correlated Poisson Sampling (Grafström, 2012), Balanced Ac-
 165 ceptance Sampling (Robertson et al., 2013), Within-Sample-Distance Sampling
 166 (Benedetti and Piersimoni, 2017), and Halton Iterative Partitioning Sampling
 167 (Robertson et al., 2018). In this manuscript, we select spatially balanced sam-
 168 ples using the Generalized Random Tessellation Stratified (GRTS) algorithm
 169 because it has several attractive properties: the GRTS algorithm accommodates
 170 finite and infinite sampling frames, equal, unequal, and proportional (to size)
 171 inclusion probabilities, legacy (historical) sampling (Foster et al., 2017), a min-
 172 imum distance between units in a sample, and replacement units (replacement
 173 units are population units that can be sampled when a population unit origi-
 174 nally selected can no longer be sampled). The GRTS algorithm selects samples
 175 by utilizing a particular mapping between two-dimensional and one-dimensional
 176 space that preserves proximity relationships. Via this mapping, units in two-
 177 dimensional space are partitioned using a hierarchical address. This hierarchical
 178 address is used to map population units to a one-dimensional line. On the one
 179 dimensional line, each population unit’s line length equals its inclusion proba-
 180 bility. Then, a systematic sample of population units is selected on the line and
 181 mapped back to two-dimensional space, yielding the desired sample. Stevens
 182 and Olsen (2004) provide more technical details.

After selecting a sample and collecting data, unbiased estimates of popu-
 lation means and totals can be obtained using the Horvitz-Thompson estima-
 tor (Horvitz and Thompson, 1952). If τ is a population total, the Horvitz-

Thompson estimator for τ , denoted by $\hat{\tau}_{ht}$, is given by

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

where Z_i is the value of the i th population unit in the sample, π_i is the inclusion probability of the i th population unit in the sample, and n is the sample size. 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 the uncertainty in $\hat{\tau}_{ht}$. 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 more precise (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 technical details, see

206 Ver Hoef (2008). Let $\mathbf{z} \equiv \{z(s_1), z(s_2), \dots, z(s_N)\}$ be an $N \times 1$ response vector
 207 at locations s_1, s_2, \dots, s_N that can be measured at the N population
 208 units. Suppose we want to use a sample to predict some linear function of the
 209 response variable, $f(\mathbf{z}) = \mathbf{b}'\mathbf{z}$, where \mathbf{b}' is a $1 \times N$ vector of weights (e.g, the
 210 population mean is represented by a weights vector whose elements all equal
 211 $1/N$). Denoting quantities that are part of the sampled population units with a
 212 subscript s and quantities that are part of the unsampled population units with
 213 a subscript u , let

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

214 where \mathbf{X}_s and \mathbf{X}_u are the design matrices for the sampled and unsampled
 215 population units, respectively, $\boldsymbol{\beta}$ is the parameter vector of fixed effects, and
 216 $\boldsymbol{\delta} \equiv [\boldsymbol{\delta}_s \ \boldsymbol{\delta}_u]'$, where $\boldsymbol{\delta}_s$ and $\boldsymbol{\delta}_u$ are random errors for the sampled and unsampled
 217 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 (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 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. Distributional assumptions are used, however, when constructing prediction intervals.

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

241 data contexts, and FBPK tended to perform best.

242 **2. Materials and Methods**

243 *2.1. Simulation Study*

244 We used a simulation study to investigate performance of four sampling-
245 analysis combinations. The first sampling-analysis combination was IRS-Design.
246 In IRS-Design, samples were selected with the Independent Random Sampling
247 (IRS) algorithm. The IRS algorithm ignores the spatial locations of the pop-
248 ulation units, thus the IRS samples were not spatially balanced. In IRS-
249 Design, samples were analyzed using the design-based approach via the Horvitz-
250 Thompson mean estimator and an IRS variance estimator that ignored the spa-
251 tial locations of the units in the sample. The second sampling-analysis combi-
252 nation was IRS-Model, where samples were selected with the IRS algorithm and
253 analyzed using the model-based approach via Restricted Maximum Likelihood
254 (REML) estimation (Harville, 1977; Patterson and Thompson, 1971; Wolfinger
255 et al., 1994). The third sampling-analysis combination was GRTS-Design,
256 where samples were selected with the GRTS algorithm and analyzed using the
257 design-based approach via the Horvitz-Thompson mean estimator and the local
258 neighborhood variance estimator (which does incorporate the spatial locations
259 of the units in the sample). The fourth and final sampling-analysis combina-
260 tion was GRTS-Model, where samples were selected with the GRTS algorithm
261 and analyzed using the model-based approach via REML estimation. These
262 sampling-analysis combinations are also provided in Table 1. Lastly we note
263 that for both the IRS and GRTS samples, equal inclusion probabilities were as-
264 sumed for all population units. When IRS assumes equal inclusion probabilities
265 for all population units, the algorithm is equivalent to simple random sampling
266 (SRS).

	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 (of the population units), 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 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) from Equation 3, respectively. The total variance, $\sigma_1^2 + \sigma_2^2$, was always 2. The range was always $\sqrt{2}/3$, chosen so that the correlation in the dependent random error decayed to nearly zero at $\sqrt{2}$, 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 was 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 95% confidence (design-based) or 95% prediction (model-based) intervals. Then we recorded the bias, squared error, standard error, and interval coverage for all sampling-analysis combinations. After all 2000 trials, we summarized the long-run performance of the combinations by calculating mean bias, rMS(P)E (root-mean-squared error for the design-based approaches and root-mean-squared-prediction error for the model-based approaches), MStdE (mean standard error), and the proportion of times the true mean is contained in its 95% confidence (design-based) or 95% prediction (model-based) interval. The 95% intervals were constructed using the normal distribution. Justification for this comes from the asymptotic normality of means via the Central Limit Theorem (under some assumptions). Quantifying mean bias and rMS(P)E is important because they help us understand how far (under different loss metrics) the estimates (design-based) or predictions (model-based) tend to be from the true mean. Quantifying MStdE is important because it helps us understand how precise intervals tend to be. Quantifying interval coverage is important because it helps us understand how often our 95% intervals actually contain the true mean.

The IRS algorithm, IRS variance estimator, GRTS algorithm, and local

neighborhood variance estimator are available in the **spsurvey R** package (Dumelle et al., 2021). FPBK is available in the **sptotal R** package (Higham et al., 2021).

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. One component of NARS is the National Lakes Assessment (NLA), which measures various aspects of lake health and water quality (USEPA, 2012). 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-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% intervals from this sample of 100 lakes and compare to the true mean mercury concentration from all 986 lakes.

3. Results

3.1. Simulation Study

The mean bias was nearly zero for all four sampling-analysis 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. 7 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

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

336 When there is no spatial covariance (Fig. 7, “Prop DE: 0” row), the four
 337 sampling-analysis combinations have approximately equal rMS(P)E and using
 338 the GRTS algorithm or a model-based analysis does not result in much, if any,
 339 loss in efficiency compared to IRS-Design. When there is spatial covariance
 340 (Fig. 7, “Prop DE: 0.5” and “Prop DE: 0.9” rows), GRTS-Model tends to
 341 have the lowest rMS(P)E, followed by GRTS-Design, IRS-Model, and finally
 342 IRS-Design, though the difference in relative rMS(P)E among GRTS-Model,
 343 GRTS-Design, and IRS-Model is relatively small. As the strength of spatial
 344 covariance increases, the gap in rMS(P)E between IRS-Design and the other
 345 sampling-analysis combinations widens. Finally we note that when there is
 346 spatial covariance, IRS-Model has a much lower rMS(P)E than IRS-Design,
 347 suggesting that the poor design properties of IRS are largely mitigated by the
 348 model-based analysis. These rMS(P)E conclusions are similar to those observed
 349 in the grid location layout, so we omit a grid location layout figure here. Ta-
 350 bles for rMS(P)E in all 36 simulation scenarios are provided in the supporting
 351 information.

Fig. 8 shows the relative MStdE of the four sampling-analysis combina-
 tions using the random location layout with “IRS-Design” as the baseline. The
 relative MStdE is defined as

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

352 Many general takeaways regarding MStdE are similar to general takeaways re-
 353 garding rMS(P)E: there seems to be no benefit to using IRS, even when there

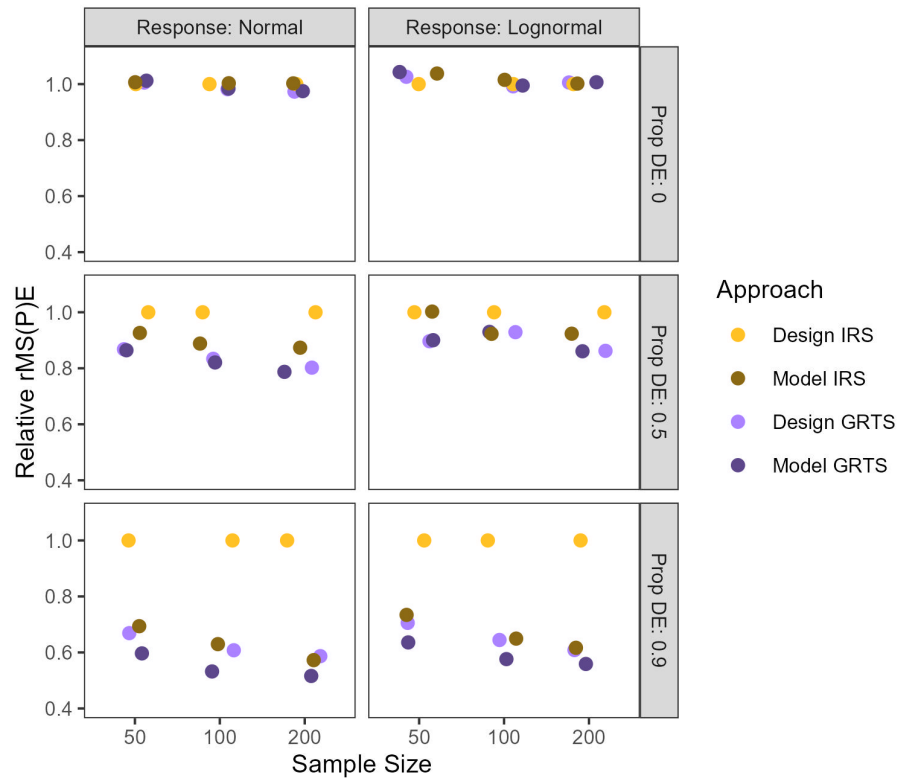


Figure 2: Relative rMS(P)E in the simulation study for the four sampling-analysis combinations. The rows indicate the proportion of dependent error and the columns indicate the response type.

is no spatial covariance; as the strength of spatial covariance increases, the gap in MStdE between IRS-Design and the other sampling-analysis 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-analysis 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. 8, “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. 8, “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. 9 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, 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-analysis 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.

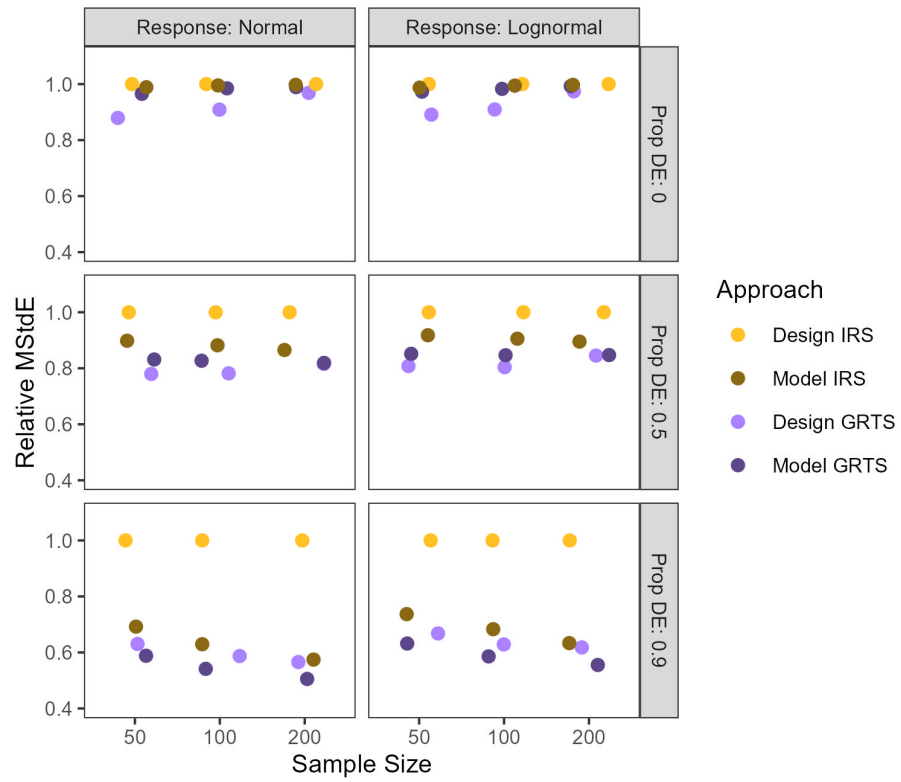


Figure 3: Relative MStdE in the simulation study for the four sampling-analysis combinations. The rows indicate the proportion of dependent error and the columns indicate the response type.

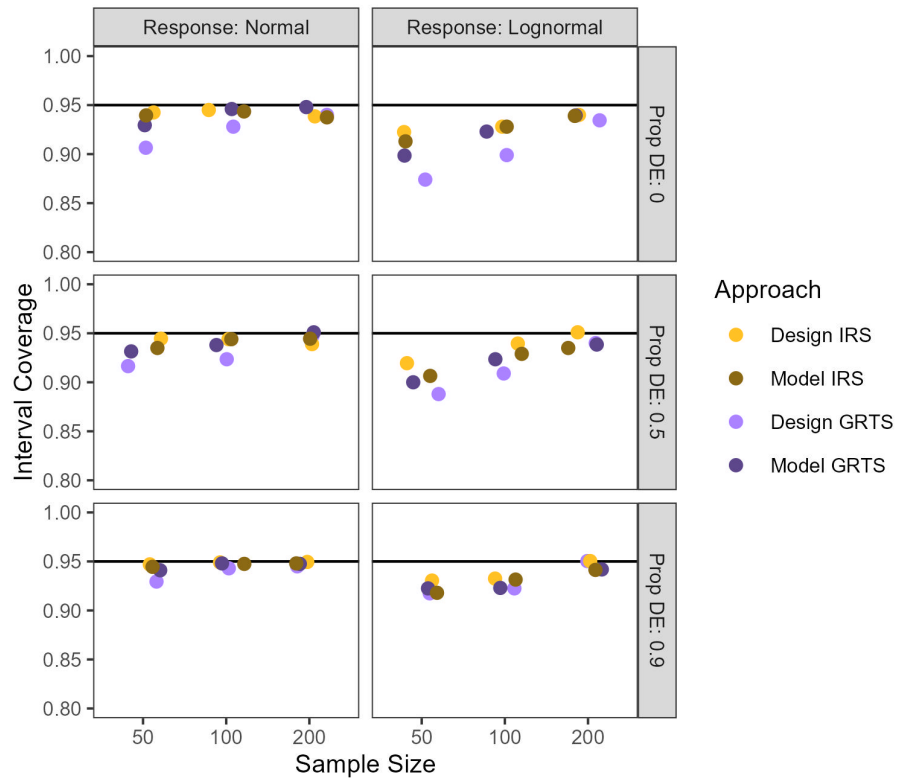


Figure 4: Interval coverage in the simulation study for the four sampling-analysis combinations. The rows indicate the proportion of dependent error and the columns indicate the response type. The solid, black line represents 95% coverage.

3.2. Application

Fig. 6 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. 6 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. 6 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 3 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-analysis combinations, the true mean mer-

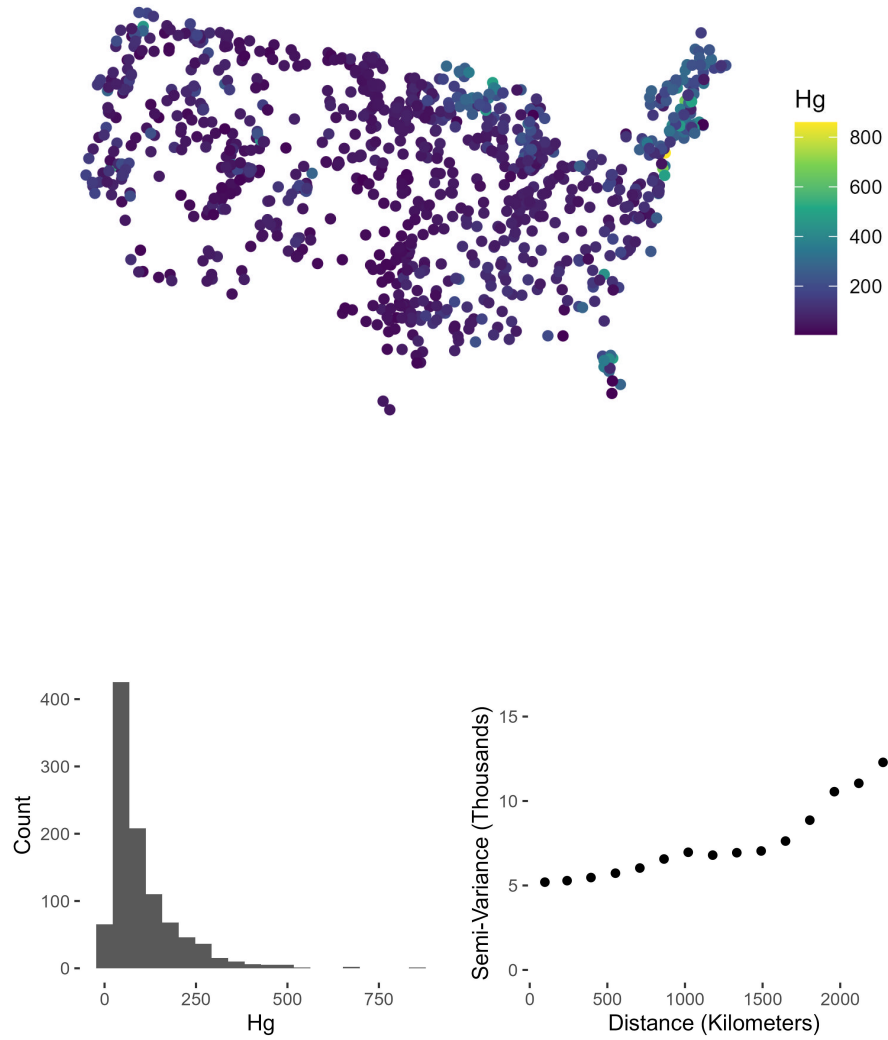


Figure 5: Mercury concentration (Hg) visualizations for all 986 lakes in the NLA data. A spatial layout is in the top row, a histogram is in the bottom row and left column, and an empirical semivariogram is in the bottom row and right column.

cury 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 3: For each sampling-analysis 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

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 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

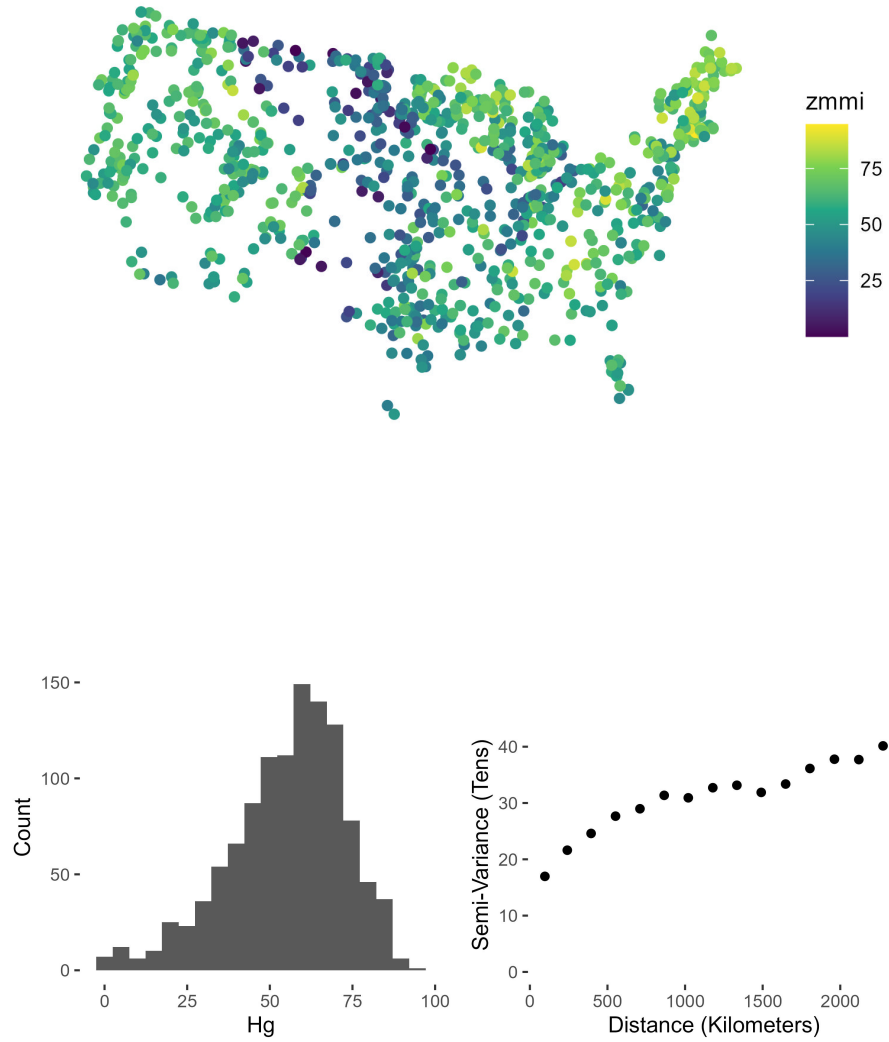


Figure 6: zmmi visualizations for all 986 lakes in the NLA data. A spatial layout is in the top row, a histogram is in the bottom row and left column, and an empirical semivariogram is in the bottom row and right column.

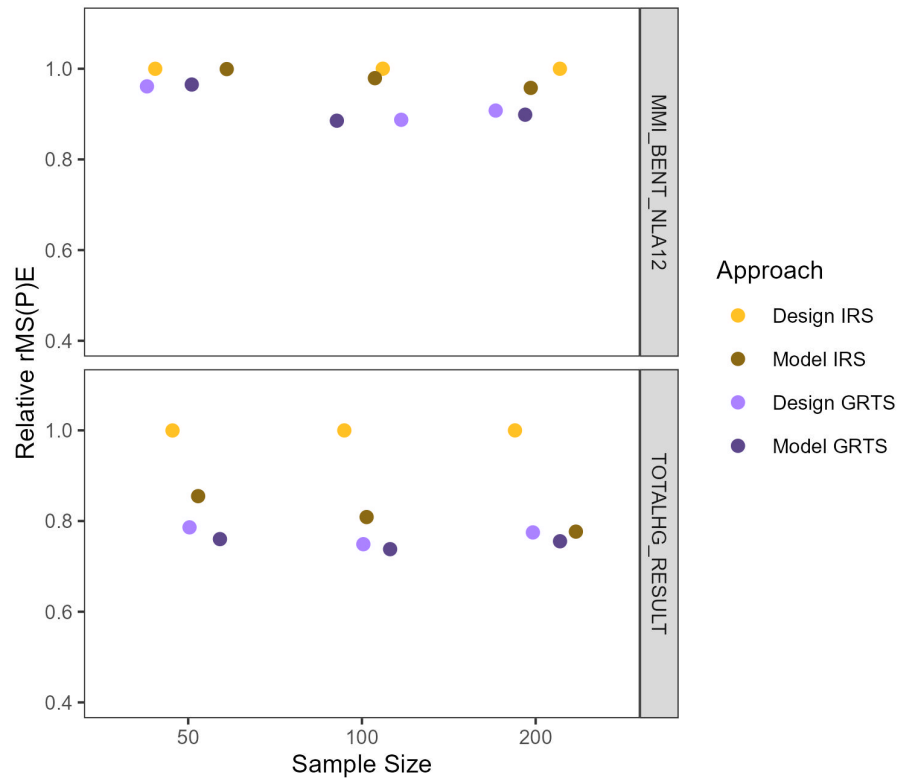


Figure 7: Relative RMS(P)E in the data study for the four sampling-analysis combinations. The rows indicate the proportion of dependent error and the columns indicate the response type.

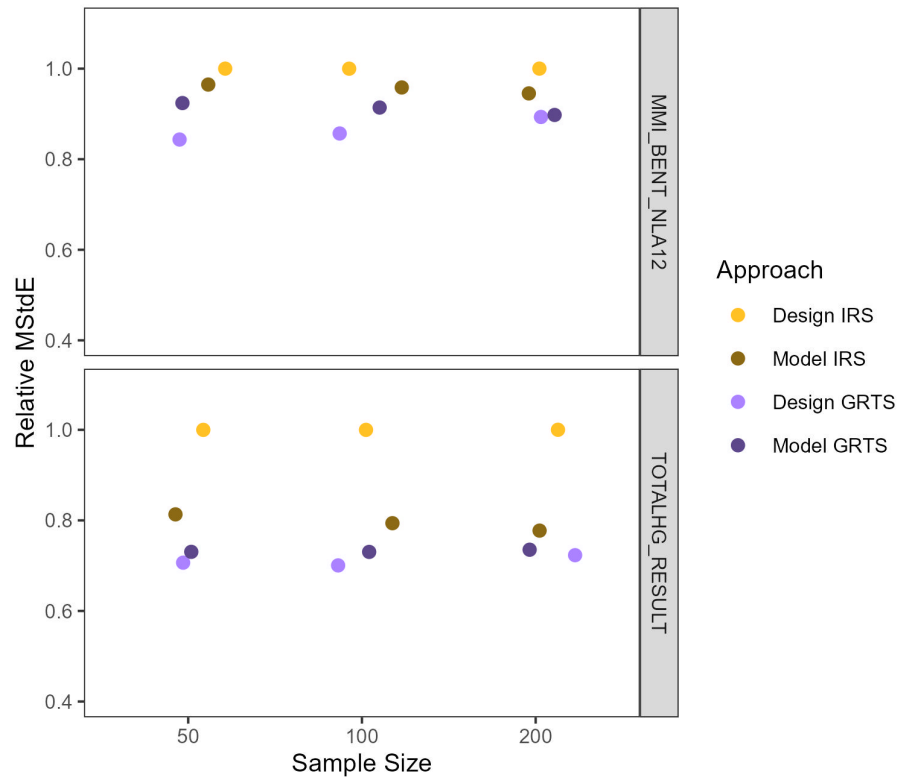


Figure 8: Relative MStdE in the data study for the four sampling-analysis combinations. The rows indicate the proportion of dependent error and the columns indicate the response type.

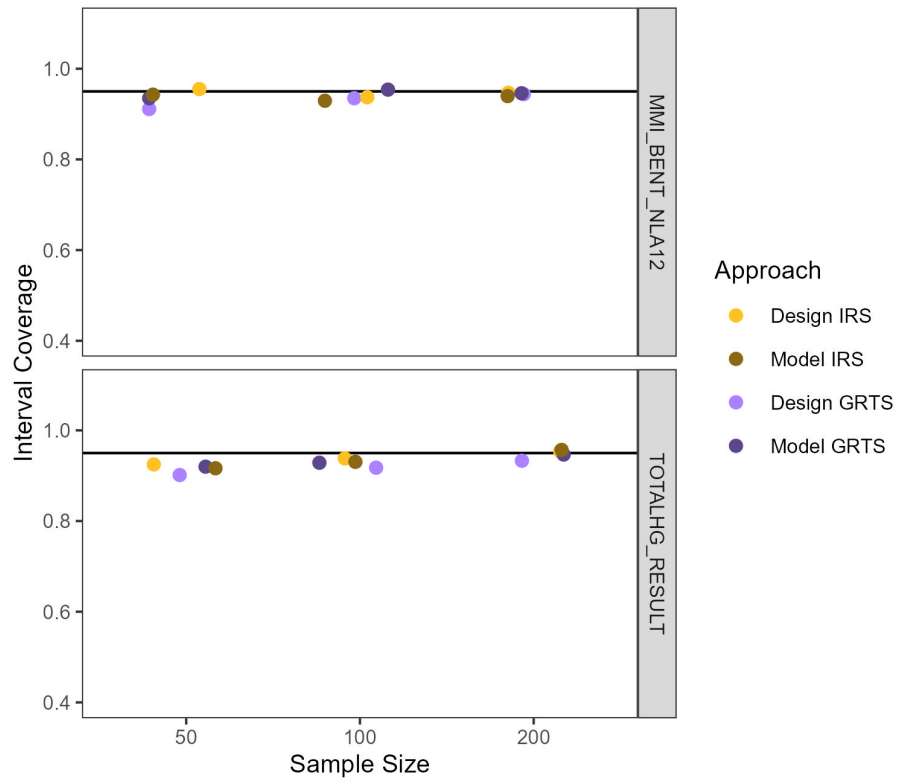


Figure 9: Interval coverage in the data study for the four sampling-analysis combinations. The rows indicate the proportion of dependent error and the columns indicate the response type. The solid, black line represents 95% coverage.

428 decision (IRS vs GRTS) is most important when using a design-based analysis.
 429 Though GRTS-Model still had lower rMS(P)E than IRS-Model, the model-based
 430 analysis mitigated most of the rMS(P)E inefficiencies that result from the IRS
 431 samples lacking spatial balance. Third, as the strength of spatial covariance
 432 increases, the gap in rMS(P)E and MStdE between IRS-Design and the other
 433 sampling-analysis combinations also increases, likely because IRS-Design is the
 434 only combination that ignores spatial locations in sampling and analysis. Fourth
 435 and finally, when the response was normal, interval coverage for all sampling-
 436 analysis combinations was usually close to 95% for all sample sizes; when the
 437 response was lognormal, interval coverage for all sampling-analysis combinations
 438 was usually between 90% and 95% and closest to 95% when $n = 200$.

439 There are several benefits and drawbacks of the design-based and model-
 440 based approaches for finite population spatial data. Some we have discussed,
 441 but others we have not, and they are worthy of consideration in future research.
 442 Design-based approaches are often computationally efficient, while model-based
 443 approaches can be computationally burdensome, especially for likelihood-based
 444 estimation methods like REML that rely on inverting a covariance matrix. The
 445 design-based approach also more naturally handles binary data, free from the
 446 more complicated logistic regression framework commonly used to analyze bi-
 447 nary data in a model-based approach. The model-based approach, however, can
 448 more naturally quantify the relationship between covariates (predictor variables)
 449 and the response variable. The model-based approach also yields estimated
 450 spatial covariance parameters, which help better understand the dependence
 451 structure in the stochastic process of study. Model selection is also possible
 452 using model-based approaches and criteria such as cross validation, likelihood
 453 ratio tests, or AIC (Akaike, 1974). Model-based approaches are capable of
 454 more efficient small-area estimation than design-based approaches by leverag-

ing distributional assumptions in areas with few observed units. Model-based approaches can also compute unit-by-unit predictions at unobserved locations and use them to construct informative visualizations like smoothed maps. In short, when deciding whether the design-based or model-based approach is more appropriate to implement, the benefits and drawbacks of each approach should be considered alongside the particular goals of the study.

Acknowledgments

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

Conflict of Interest Statement

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

Author Contribution Statement

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

477 **Data and Code Availability**

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

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

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

483 **Supporting Information**

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

486 **References**

- 487 Akaike, H., 1974. A new look at the statistical model identification. IEEE
 488 Transactions on Automatic Control 19, 716–723.
- 489 Barabesi, L., Franceschi, S., 2011. Sampling properties of spatial total estima-
 490 tors under tessellation stratified designs. Environmetrics 22, 271–278.
- 491 Benedetti, R., Piersimoni, F., 2017. A spatially balanced design with probability
 492 function proportional to the within sample distance. Biometrical Journal 59,
 493 1067–1084.
- 494 Benedetti, R., Piersimoni, F., Postiglione, P., 2017. Spatially balanced sam-
 495 pling: A review and a reappraisal. International Statistical Review 85, 439–
 496 454.
- 497 Breiman, L., 2001. Random forests. Machine Learning 45, 5–32.
- 498 Brus, D., De Gruijter, J., 1997. Random sampling or geostatistical modelling?
 499 Choosing between design-based and model-based sampling strategies for soil
 500 (with discussion). Geoderma 80, 1–44.

- 501 Brus, D.J., 2021. Statistical approaches for spatial sample survey: Persistent
502 misconceptions and new developments. *European Journal of Soil Science* 72,
503 686–703.
- 504 Chan-Golston, A.M., Banerjee, S., Handcock, M.S., 2020. Bayesian inference for
505 finite populations under spatial process settings. *Environmetrics* 31, e2606.
- 506 Chiles, J.-P., Delfiner, P., 1999. *Geostatistics: Modeling Spatial Uncertainty*.
507 John Wiley & Sons, New York.
- 508 Cicchitelli, G., Montanari, G.E., 2012. Model-assisted estimation of a spatial
509 population mean. *International Statistical Review* 80, 111–126.
- 510 Cooper, C., 2006. Sampling and variance estimation on continuous domains.
511 *Environmetrics* 17, 539–553.
- 512 Cressie, N., 1993. *Statistics for spatial data*. John Wiley & Sons.
- 513 De Gruijter, J., Ter Braak, C., 1990. Model-free estimation from spatial sam-
514 ples: A reappraisal of classical sampling theory. *Mathematical Geology* 22,
515 407–415.
- 516 Diggle, P.J., Menezes, R., Su, T., 2010. Geostatistical inference under prefer-
517 ential sampling. *Journal of the Royal Statistical Society: Series C (Applied*
518 *Statistics)* 59, 191–232.
- 519 Dumelle, M., Kincaid, T.M., Olsen, A.R., Weber, M.H., 2021. *Spsurvey: Spatial*
520 *sampling design and analysis*.
- 521 Fix, E., Hodges, J.L., 1989. Discriminatory analysis. Nonparametric discrimi-
522 nation: Consistency properties. *International Statistical Review/Revue In-*
523 *ternationale de Statistique* 57, 238–247.
- 524 Foster, S.D., Hosack, G.R., Lawrence, E., Przeslawski, R., Hedge, P., Caley,
525 M.J., Barrett, N.S., Williams, A., Li, J., Lynch, T., others, 2017. Spatially
526 balanced designs that incorporate legacy sites. *Methods in Ecology and*
527 *Evolution* 8, 1433–1442.

- 528 Grafström, A., 2012. Spatially correlated poisson sampling. *Journal of Statis-*
529 *tistical Planning and Inference* 142, 139–147.
- 530 Grafström, A., Lundström, N.L., 2013. Why well spread probability samples
531 are balanced. *Open Journal of Statistics* 3, 36–41.
- 532 Grafström, A., Lundström, N.L., Schelin, L., 2012. Spatially balanced sampling
533 through the pivotal method. *Biometrics* 68, 514–520.
- 534 Grafström, A., Matei, A., 2018. Spatially balanced sampling of continuous
535 populations. *Scandinavian Journal of Statistics* 45, 792–805.
- 536 Hansen, M.H., Madow, W.G., Tepping, B.J., 1983. An evaluation of model-
537 dependent and probability-sampling inferences in sample surveys. *Journal*
538 *of the American Statistical Association* 78, 776–793.
- 539 Harville, D.A., 1977. Maximum likelihood approaches to variance component
540 estimation and to related problems. *Journal of the American Statistical*
541 *Association* 72, 320–338.
- 542 Higham, M., Ver Hoef, J., Frank, B., Dumelle, M., 2021. Sptotal: Predicting
543 totals and weighted sums from spatial data.
- 544 Horvitz, D.G., Thompson, D.J., 1952. A generalization of sampling without
545 replacement from a finite universe. *Journal of the American Statistical As-*
546 *sociation* 47, 663–685.
- 547 Lohr, S.L., 2009. *Sampling: Design and analysis*. Nelson Education.
- 548 Patterson, H.D., Thompson, R., 1971. Recovery of inter-block information when
549 block sizes are unequal. *Biometrika* 58, 545–554.
- 550 Robertson, B., Brown, J., McDonald, T., Jaksons, P., 2013. BAS: Balanced
551 acceptance sampling of natural resources. *Biometrics* 69, 776–784.
- 552 Robertson, B., McDonald, T., Price, C., Brown, J., 2018. Halton iterative
553 partitioning: Spatially balanced sampling via partitioning. *Environmental*
554 *and Ecological Statistics* 25, 305–323.

- 555 Särndal, C.-E., Swensson, B., Wretman, J., 2003. Model assisted survey sam-
556 pling. Springer Science & Business Media.
- 557 Schabenberger, O., Gotway, C.A., 2017. Statistical methods for spatial data
558 analysis. CRC press.
- 559 Sen, A.R., 1953. On the estimate of the variance in sampling with varying
560 probabilities. *Journal of the Indian Society of Agricultural Statistics* 5, 127.
- 561 Sterba, S.K., 2009. Alternative model-based and design-based frameworks for
562 inference from samples to populations: From polarization to integration.
563 *Multivariate Behavioral Research* 44, 711–740.
- 564 Stevens, D.L., Olsen, A.R., 2003. Variance estimation for spatially balanced
565 samples of environmental resources. *Environmetrics* 14, 593–610.
- 566 Stevens, D.L., Olsen, A.R., 2004. Spatially balanced sampling of natural re-
567 sources. *Journal of the American Statistical Association* 99, 262–278.
- 568 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)
569 [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)
570 [assessment](https://www.epa.gov/national-aquatic-resource-surveys/national-results-and-regional-highlights-national-lakes-assessment).
- 571 Ver Hoef, J., 2002. Sampling and geostatistics for spatial data. *Ecoscience* 9,
572 152–161.
- 573 Ver Hoef, J.M., 2008. Spatial methods for plot-based sampling of wildlife pop-
574 ulations. *Environmental and Ecological Statistics* 15, 3–13.
- 575 Ver Hoef, J.M., Temesgen, H., 2013. A comparison of the spatial linear model
576 to nearest neighbor (k-NN) methods for forestry applications. *PIOS ONE* 8,
577 e59129.
- 578 Wang, J.-F., Jiang, C.-S., Hu, M.-G., Cao, Z.-D., Guo, Y.-S., Li, L.-F., Liu,
579 T.-J., Meng, B., 2013. Design-based spatial sampling: Theory and imple-
580 mentation. *Environmental Modelling & Software* 40, 280–288.

- 581 Wang, J.-F., Stein, A., Gao, B.-B., Ge, Y., 2012. A review of spatial sampling.
582 Spatial Statistics 2, 1–14.
- 583 Wolfinger, R., Tobias, R., Sall, J., 1994. Computing gaussian likelihoods and
584 their derivatives for general linear mixed models. SIAM Journal on Scientific
585 Computing 15, 1294–1310.