

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.
2.
3.
4. 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. In this manuscript, we compare the approaches for finite population spatial data. We first provide relevant background for the approaches and then use a simulation study and an analysis of real mercury concentration data to numerically compare them. We find that sampling plans that incorporate spatial locations (spatially balanced samples) perform better than sampling plans ignoring spatial locations (non-spatially balanced samples), regardless of whether design-based or model-based approaches were used to analyze the data. We also find that within sampling plans, the model-based approaches often outperform design-based approaches, even for skewed data. This gap in performance is small when spatially balanced samples

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

are used but large when non-spatially balanced samples are used.

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 is called a sample. In the design-based approach, inferences about the underlying population are informed via a probabilistic process assigning some population units to be part of 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;

55 Wang et al., 2012). Cooper (2006) reviews the two approaches in an ecological
56 context before introducing a “model-assisted” variance estimator that combines
57 aspects from each approach. In addition to Cooper (2006), there has been
58 substantial research and development into estimators that use both design and
59 model-based principles (see e.g., Sterba (2009), Cicchitelli and Montanari (2012),
60 Chan-Golston et al. (2020) for a Bayesian approach).

61 Certainly comparisons between design-based and model-based approaches to
62 spatial data have been studied. But no numerical comparison has been made
63 between design-based approaches incorporating spatial information and design-
64 based approaches. In this manuscript, we compare design-based approaches
65 incorporating spatial information to model-based approaches for spatial data.
66 We focus on finite populations, but these comparisons generalize to infinite
67 populations as well. A finite population contains a finite number of population
68 units; an example is lakes (treated as a whole with the lake centroid representing
69 location) in the contiguous United States. An infinite population contains an
70 infinite number of population units; an example is locations within a single lake.

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

79 1.1. Background

80 The design-based and model-based approaches incorporate randomness in
81 fundamentally different ways. In this section, we describe the role of randomness

82 for each approach and the subsequent effects on statistical inferences for spatial
83 data.

84 *1.1.1. Comparing Design-Based and Model-Based Approaches*

85 The design-based approach assumes the population is fixed. Randomness is
86 incorporated via the selection of units in a sampling frame. A sampling frame
87 is the set of all units available to be sampled. Units from the sampling frame
88 are selected as part of the sample according to a sampling design, which assigns
89 a positive probability of inclusion (inclusion probability) to each unit from the
90 sampling frame. Some examples of commonly used sampling designs include
91 simple random sampling, stratified random sampling, and cluster sampling.
92 When sampling designs incorporate spatial locations into sampling, we call
93 the resulting samples “spatially balanced.” One approach to selecting spatially
94 balanced samples is the Generalized Random Tessellation Stratified (GRTS)
95 algorithm (Stevens and Olsen, 2004), which we discuss in more detail in Section
96 1.1.2. When sampling designs do not incorporate spatial locations into sampling,
97 we call the resulting samples “non-spatially balanced.”

98 Fundamentally, the design-based approach combines the randomness of the
99 sampling design with the data collected via the sample to justify the estimation
100 and uncertainty quantification of fixed, unknown parameters of a population (e.g.,
101 a population mean). Treating the data as fixed and incorporating randomness
102 through the sampling design yields estimators having very few other assumptions.
103 Confidence intervals for these types of estimators are typically derived using
104 limiting arguments that incorporate all possible samples. Sample means, for
105 example, are asymptotically normal (Gaussian) by the Central Limit Theorem
106 (under some assumptions). If we repeatedly select samples from the population,
107 then 95% of all 95% confidence intervals constructed from a procedure with
108 appropriate coverage will contain the true, fixed mean. Särndal et al. (2003)

109 and Lohr (2009) provide thorough reviews of the design-based approach.

110 The model-based approach assumes the data are a random realization of
111 a data-generating stochastic process. Randomness is incorporated through
112 distributional assumptions on this process. Strictly speaking, randomness need
113 not be incorporated through random sampling, though Diggle et al. (2010) warn
114 against preferential sampling. Preferential sampling occurs when the process
115 generating the data locations and the process being modeled are not independent
116 of one another. To guard against preferential sampling, model-based approaches
117 often still implement some form of random sampling.

118 Instead of estimating fixed, unknown population parameters, as in the design-
119 based approach, often the goal of model-based inference is to predict a realized
120 variable, or value. For example, suppose the realized mean of all population
121 units is the value of interest. Instead of *estimating* a fixed, unknown mean, we
122 are *predicting* the value of the mean, a random variable. Prediction intervals are
123 then derived using assumptions of the data-generating stochastic process. If we
124 repeatedly generate response values from the same data-generating stochastic
125 process and select samples, then 95% of all 95% prediction intervals constructed
126 from a procedure with appropriate coverage will contain their respective realized
127 means. Cressie (1993) and Schabenberger and Gotway (2017) provide thorough
128 reviews of model-based approaches for spatial data. In Fig. 1, we provide a
129 visual comparison of the design-based and model-based approaches (Ver Hoef
130 (2002) and Brus (2021) provide similar figures).

131 1.1.2. Spatially Balanced Design and Analysis

132 We previously mentioned that the design-based approach can be used to
133 select spatially balanced samples (samples that incorporate spatial locations
134 of the population units and are “well-spread” in space). Spatially balanced
135 samples are useful because parameter estimates from these samples tend to

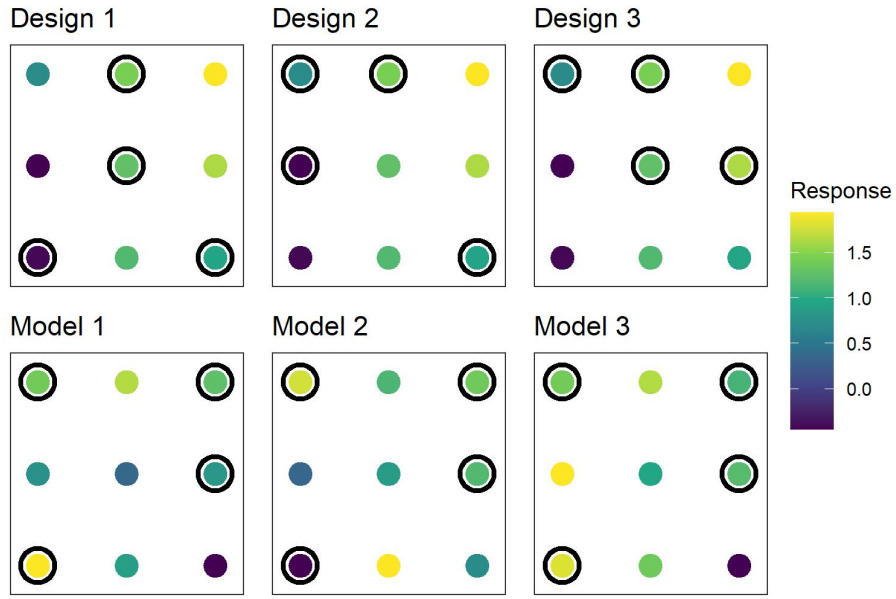


Figure 1: A visual comparison of the design-based and model-based approaches. In the top row, there is one fixed population with nine population units and three random samples of size four (points circled are those sampled). The response values at each site are fixed, but we obtain different estimates 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

136 vary less than parameter estimates from samples that are not spatially balanced
 137 (Barabesi and Franceschi, 2011; Benedetti et al., 2017; Grafström and Lundström,
 138 2013; Robertson et al., 2013; Stevens and Olsen, 2004; Wang et al., 2013).
 139 The first spatially balanced sampling algorithm seeing widespread use is the
 140 Generalized Random Tessellation Stratified (GRTS) algorithm (Stevens and
 141 Olsen, 2004). To quantify the spatial balance of a sample, Stevens and Olsen
 142 (2004) proposed loss metrics based on Voronoi polygons (Dirichlet Tessellations).
 143 After the GRTS algorithm was developed, several other spatially balanced
 144 sampling algorithms emerged, such as the Local Pivotal Method (Grafström et
 145 al., 2012; Grafström and Matei, 2018), Spatially Correlated Poisson Sampling
 146 (Grafström, 2012), Balanced Acceptance Sampling (Robertson et al., 2013),
 147 Within-Sample-Distance Sampling (Benedetti and Piersimoni, 2017), and Halton
 148 Iterative Partitioning Sampling (Robertson et al., 2018). In this manuscript, we
 149 select spatially balanced samples using the Generalized Random Tessellation
 150 Stratified (GRTS) algorithm because it has several attractive properties. More
 151 specifically, the GRTS algorithm accommodates finite and infinite sampling
 152 frames, equal, unequal, and proportional (to size) inclusion probabilities, legacy
 153 (historical) sampling (Foster et al., 2017), a minimum distance between units in
 154 a sample, and replacement units (replacement units are population units that
 155 can be sampled when a population unit originally selected can no longer be
 156 sampled). The GRTS algorithm selects samples by utilizing a particular mapping
 157 between two-dimensional and one-dimensional space that preserves proximity
 158 relationships. Via this mapping, units in two-dimensional space are partitioned
 159 using a hierarchical address. This hierarchical address is used to map population
 160 units to a one-dimensional line. On the one dimensional line, each population
 161 unit's line length equals its inclusion probability. Then, a systematic sample of
 162 population units is selected on the line, yielding desired sample. Stevens and

163 Olsen (2004) provides 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 estimate of τ , denoted by $\hat{\tau}_{ht}$, is given by

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

164 where Z_i is the value of the i th population unit in the sample and π_i is the
 165 inclusion probability of the i th population unit in the sample. An estimate of
 166 the population mean is obtained by dividing $\hat{\tau}_{ht}$ by N , the number of population
 167 units.

168 It is also important to quantify uncertainty $\hat{\tau}_{ht}$. Horvitz and Thompson
 169 (1952) and Sen (1953) provide variance estimators for $\hat{\tau}_{ht}$, but these estimators
 170 have two drawbacks. First, they rely on calculating π_{ij} , the probability that
 171 population unit i and population unit j are both in the sample – this quantity
 172 can be challenging if not impossible to calculate analytically. Second, these
 173 estimators ignore the spatial locations of the population units. To address these
 174 two drawbacks simultaneously, Stevens and Olsen (2003) proposed the local
 175 neighborhood variance estimator. The local neighborhood variance estimator
 176 does not rely on π_{ij} and incorporates spatial locations – for technical details see
 177 Stevens and Olsen (2003). Stevens and Olsen (2003) show the local neighborhood
 178 variance estimator tends to reduce the estimated variance of $\hat{\tau}$ and yield narrower
 179 confidence intervals compared to variance estimators that ignore spatial locations.

180 1.1.3. Finite Population Block Kriging

181 Finite Population Block Kriging (FPBK) is a model-based approach that
 182 expands the geostatistical Kriging framework to the finite population setting

183 (Ver Hoef, 2008). Instead of developing inference based on a specific sampling
 184 design, we assume the data are generated by a spatial stochastic process. We
 185 summarize some of the basic principles of FBPK next (for more technical details,
 186 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
 187 at locations s_1, s_2, \dots, s_N that can be measured at the N population units.
 188 Suppose we want to use a sample to predict some linear function of the response
 189 variable, $f(\mathbf{z}) = \mathbf{b}'\mathbf{z}$, where \mathbf{b}' is a $1 \times N$ vector of weights (e.g, the population
 190 mean is represented by a weights vector whose elements all equal one). Denoting
 191 quantities that are part of the sampled population units with a subscript s and
 192 quantities that are part of the unsampled population units with subscript u , let

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

193 where \mathbf{X}_s and \mathbf{X}_u are the design matrices for the sampled and unsampled
 194 population units, respectively, $\boldsymbol{\beta}$ is the parameter vector of fixed effects, and
 195 $\boldsymbol{\delta} \equiv [\boldsymbol{\delta}_s \ \boldsymbol{\delta}_u]'$, where $\boldsymbol{\delta}_s$ and $\boldsymbol{\delta}_u$ are random errors for the sampled and unsampled
 196 population units, respectively.

FBPK assumes $\boldsymbol{\delta}$ in Equation 2 has mean-zero and a spatial correlation
 structure that can be modeled using a covariance function. This covariance
 function is commonly assumed to be non-negative (between zero and one), 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 both.
 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), random forests (Breiman, 2001), Bayesian models (Chan-Golston et al., 2020), 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

219 compared FBPK, k-nearest-neighbors, and random forests in a variety of spatial
 220 data contexts, and FBPK tended to perform best. Compared to the Bayesian
 221 approach, we prefer FBPK mostly because it is much more computationally
 222 efficient.

223 2. Materials and Methods

224 2.1. Simulation Study

225 We used a simulation study to investigate performance of four sampling-
 226 analysis combinations: IRS with a design-based analysis, called “IRS-Design”;
 227 IRS with a model-based analysis, called “IRS-Model”; GRTS sampling with a
 228 design-based analysis, called “GRTS-Design”; GRTS sampling with a model-
 229 based analysis, called “GRTS-Model”. These combinations are also provided in
 230 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.

231 Performance for the four sampling-analysis combinations was evaluated in 36
 232 different simulation scenarios. The 36 scenarios resulted from the crossing of three
 233 sample sizes, two location layouts, two response types, and three proportions of
 234 dependent random error. The three sample sizes (n) were $n = 50$, $n = 100$, and
 235 $n = 200$. Samples were always selected from a population size (N) of $N = 900$.
 236 The two location layouts (of the population units) were random and gridded.
 237 Locations in the random layout were randomly generated inside the unit square
 238 $([0, 1] \times [0, 1])$. Locations in the gridded layout were placed on a fixed, equally
 239 spaced grid inside the unit square. The two response types were normal and

240 lognormal. For the normal response type, the response was simulated using mean-
 241 zero random errors with the exponential covariance (Equation 3) for varying
 242 proportions of dependent random error. The proportion of dependent random
 243 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
 244 error variance (partial sill) and independent random error variance (nugget),
 245 respectively, from Equation 3. The total variance, $\sigma_1^2 + \sigma_2^2$, was always 2. The
 246 range was always $\sqrt{2}/3$, which means that the correlation in the dependent
 247 random error decayed to nearly zero at the largest possible distance between
 248 two units in the domain. For the lognormal response type, the response was first
 249 simulated using the same approach as for the normal response type, except that
 250 the total variance was 0.6931 instead of 2. The response was then exponentiated,
 251 yielding a random variable whose total variance is 2. The lognormal responses
 252 were used to evaluate performance of the sampling-analysis approaches for data
 253 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 two.

254 In each of the 36 simulation scenarios, there were 2000 independent simulation
 255 trials. In each trial, IRS and GRTS samples were selected and then design-based
 256 and model-based analyses were used to estimate (design-based) or predict (model-
 257 based) the mean and construct confidence (design-based) or prediction (model-
 258 based) intervals. Then we recorded the bias, squared error, and interval coverage
 259 for all sampling-analysis combinations. After all 2000 trials, we summarized the
 260 long-run performance of the combinations by calculating average bias, RMS(P)E
 261 (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 Environmental Protection Agency (EPA), states, and tribes periodically conduct National Aquatic Research Surveys (NARS) in the United States to assess the water quality of various bodies of water. We will use data from the 2012 National Lakes Assessment (NLA), which measures various aspects of lake health and water quality for lakes in the contiguous United States (USEPA, 2012). Specifically, we will analyze mercury concentration in lakes. Although we know 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

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.

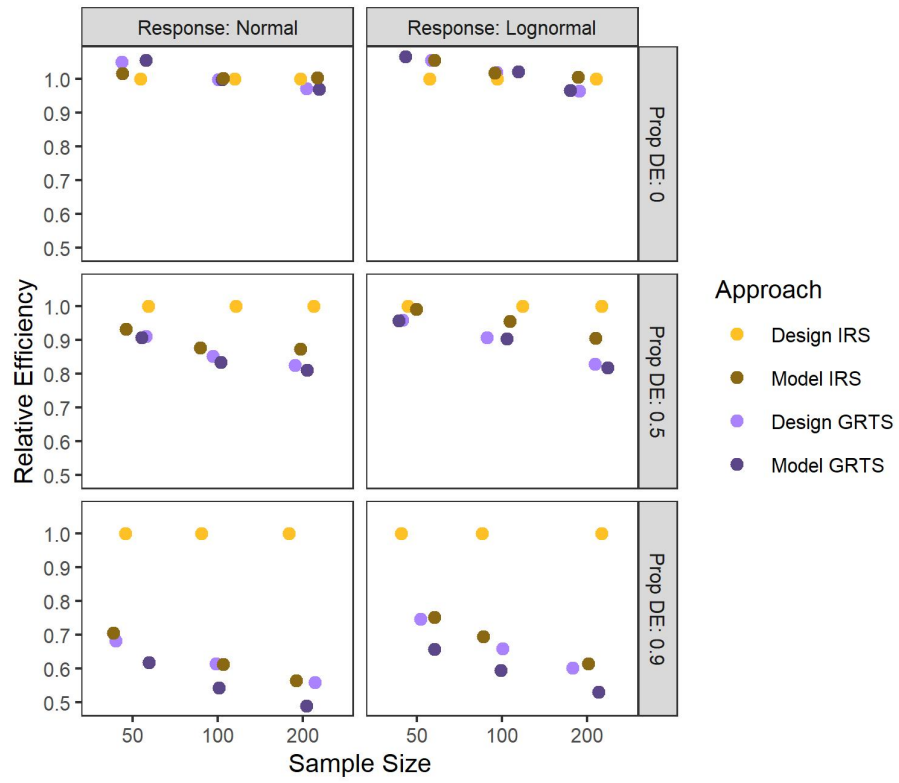


Figure 2: Relative rMS(P)E for the four sampling-analysis combinations. The rows indicate the proportion of dependent error and the columns indicate the response type.

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

We also studied 95% interval coverage among the sampling-analysis combinations. 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

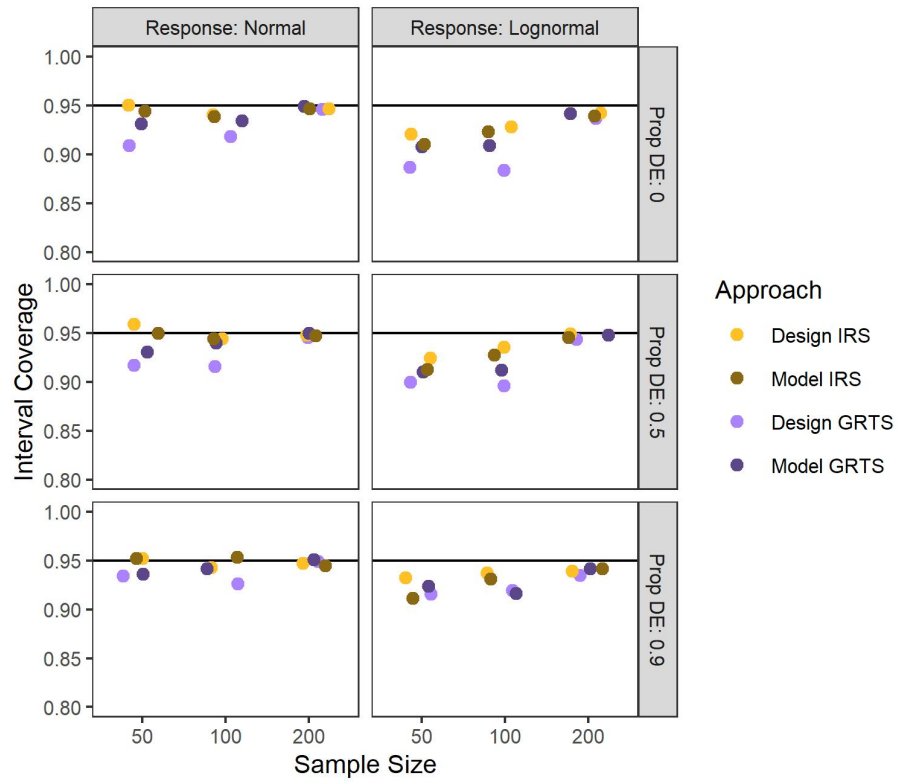


Figure 3: Interval coverage for the four sampling-analysis combinations. The rows indicate the proportion of dependent error and the columns indicate the response type. The solid, black line in each plot represents 95% coverage.

312 from from 90% to 95%. Coverage tended to always increase with the sample
313 size. At a sample size of 200, all four sampling-analysis combinations had
314 approximately 95% interval coverage in both response scenarios for all dependent
315 error proportions. These conclusions are similar to those observed in the grid
316 location layout, so we omit a grid location layout figure here. Tables for interval
317 coverage in all 36 simulation scenarios are provided in the supporting information.

318 *3.2. Application*

319 Fig. 4 shows that mercury concentration is right-skewed, with most lakes
320 having a low value of mercury concentration but a few having a much higher
321 concentration. Mercury concentration exhibits some spatial patterning, with
322 high mercury concentrations in lakes in the northeast and north central United
323 States. Fig. 4 also shows the spatial dependence in mercury concentration via
324 the empirical semivariogram. The empirical semivariogram can be used as a
325 tool to visualize spatial dependence. It quantifies the halved squared differences
326 (semivariance) among mercury concentration at different distances apart. When
327 a process is spatially correlated, the semivariance tends to be smaller at small
328 distances and larger at large distances. Together, the map, histogram, and
329 semivariogram in Fig. 4 suggest that mercury concentration is skewed and
330 exhibits spatial dependence. Lastly we note that the realized mean mercury
331 concentration in the 986 lakes is 103.2 ng / g.

332 We selected a single IRS sample and a single GRTS sample and estimated
333 (design-based) or predicted (model-based) the mean mercury concentration and
334 its standard error using using design-based and model-based approaches. For the
335 model-based analyses, the exponential covariance was used. Table 3 shows the
336 results from these analyses. For all four sampling-analysis combinations, the true
337 realized mean mercury concentration is within the bounds of the 95% confidence
338 (design-based) or prediction (model-based) intervals. Though we should not

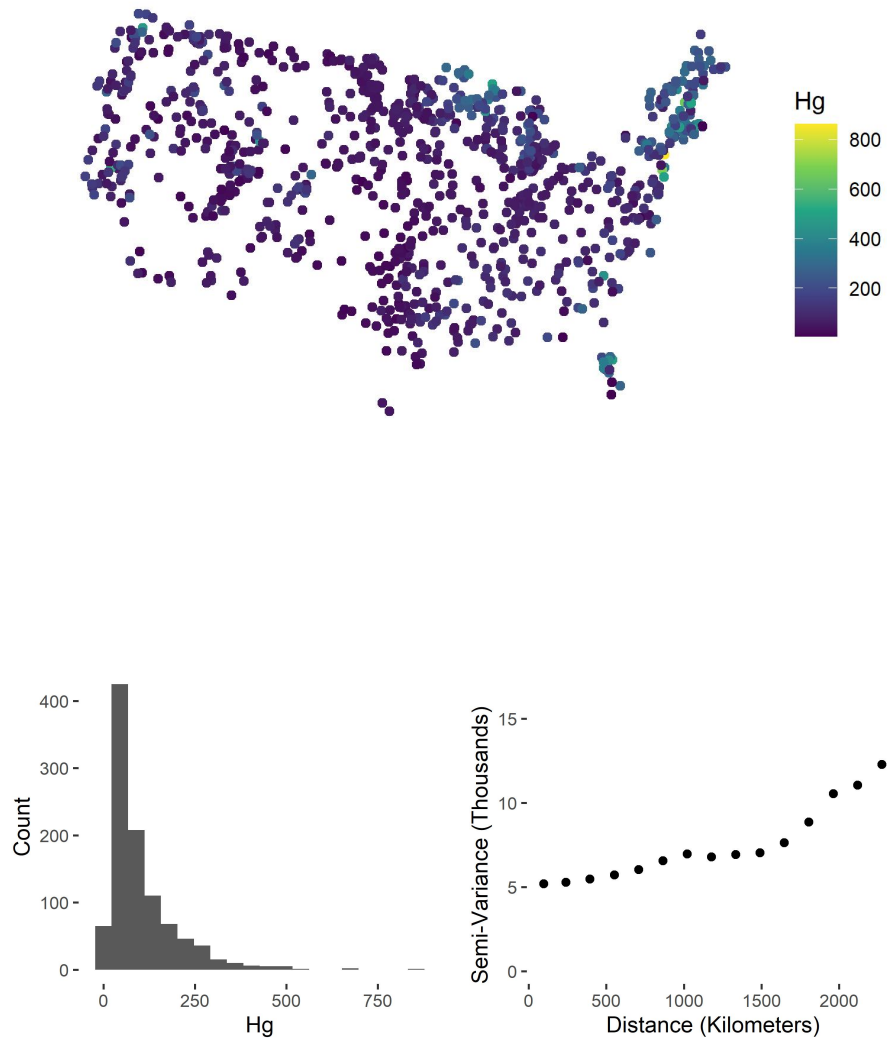


Figure 4: Mercury concentration visualizations for the population (Hg) for 986 lakes in the contiguous United States. A spatial layout is in the top row, a histogram is in the bottom row and left column, and an empirical semivariogram is in the bottom row and right column.

339 generalize these results to other samples from these data, we do note a couple
 340 of patterns. The design-based IRS analysis shows the largest standard error:
 341 a likely reason is that this is the only approach that does not incorporate any
 342 spatial information regarding mercury concentration. Both analyses using GRTS
 343 sampling have lower standard errors than both analyses using IRS sampling.
 344 We expect that these patterns are consistent with other samples from these
 345 data because mercury concentration exhibits spatial patterning, so a spatially
 346 balanced sample should usually yield a lower standard error.

Approach	Estimate	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: Application of design-based and model-based approaches to the NLA data set on
 mercury concentration. The true mean concentration is 103.2 ng / g.

347 **4. Discussion**

348 The design-based and model-based approaches to statistical inference are
 349 fundamentally different paradigms by which samples are selected and data are
 350 analyzed. The design-based approach incorporates randomness through sampling
 351 to estimate population parameters. The model-based approach incorporates
 352 randomness through distributional assumptions to predict realized values of a
 353 random process. Though these approaches have often been compared in the
 354 literature both from theoretical and analytical perspectives, our contribution
 355 lies in studying them in a spatial context while implementing spatially balanced
 356 sampling. Aside from the theoretical differences described, a few analytical
 357 findings from the simulation study are particularly notable. First, the sampling
 358 decision (GRTS vs IRS) is most important when using a design-based analysis.
 359 Though GRTS-Model still outperformed IRS-Model, the model-based analysis

360 mitigated much of the inefficiency of the IRS sample. Second, independent of
 361 the analysis approach, we found no reason to prefer IRS over GRTS for sampling
 362 spatial data – GRTS-Design and GRTS-Model generally performed at least
 363 as well as their IRS counterparts when there was no spatial correlation and
 364 noticeably better than their IRS counterparts when there was spatial correlation.
 365 Third, as the strength of spatial correlation increases, the gap in rMS(P)E
 366 between IRS-Design and the other sampling-analysis combinations also increases.
 367 Fourth and finally, when the response was normal, interval coverage for all
 368 sampling-analysis combinations was very close to 95% for all sample sizes; when
 369 the response was lognormal, interval coverage for all sampling and analysis was
 370 between 90% and 95% and closest to 95% when $n = 200$.

371 There are several benefits and drawbacks of the design-based and model-
 372 based approaches for finite population spatial data. Some we have discuss, but
 373 others we have not and they are worthy of consideration in future research.
 374 Design-based approaches are often computationally efficient, while model-based
 375 estimation can be computationally burdensome, especially for likelihood-based
 376 methods such as REML that rely on inverting a covariance matrix. The design-
 377 based approach also more naturally handles binary data, free from the more
 378 complicated logistic regression framework commonly used to analyze binary
 379 data in a model-based approach. The model-based approach, however, can
 380 more naturally quantify the relationship between covariates (predictor variables)
 381 and response variable. The model-based approach also yields estimated spatial
 382 covariance parameters, which help better understand the dependence structure
 383 in the process of study. Model selection is also possible using model-based
 384 approaches and criteria such as cross validation, likelihood ratio tests, or AIC
 385 (Akaike, 1974). Model-based approaches are capable of more efficient small-area
 386 estimation than design-based approaches by leveraging distributional assumptions

387 in areas with few observed sites. Model-based approaches can also compute site-
388 by-site predictions at unobserved locations and use them to construct informative
389 visualizations. The benefits and drawbacks of both approaches, alongside our
390 theoretical and analytical comparisons, can motivate the process of choosing among
391 them. This is especially true from an analysis perspective, as we found that
392 using a spatially balanced sampling algorithm benefits both design-based and
393 model-based analyses.

394 **Acknowledgments**

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

403 **Conflict of Interest Statement**

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

405 **Data and Code Availability**

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

410 **Supporting Information**

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

413 **Author Contributions**

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

418 **References**

- 419 Akaike, H., 1974. A new look at the statistical model identification. IEEE
420 Transactions on Automatic Control 19, 716–723.
- 421 Barabesi, L., Franceschi, S., 2011. Sampling properties of spatial total
422 estimators under tessellation stratified designs. Environmetrics 22, 271–278.
- 423 Benedetti, R., Piersimoni, F., 2017. A spatially balanced design with proba-
424 bility function proportional to the within sample distance. Biometrical Journal
425 59, 1067–1084.
- 426 Benedetti, R., Piersimoni, F., Postiglione, P., 2017. Spatially balanced
427 sampling: A review and a reappraisal. International Statistical Review 85,
428 439–454.
- 429 Breiman, L., 2001. Random forests. Machine Learning 45, 5–32.
- 430 Brus, D., De Gruijter, J., 1997. Random sampling or geostatistical modelling?
431 Choosing between design-based and model-based sampling strategies for soil
432 (with discussion). Geoderma 80, 1–44.

433 Brus, D.J., 2021. Statistical approaches for spatial sample survey: Persistent
434 misconceptions and new developments. *European Journal of Soil Science* 72,
435 686–703.

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

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

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

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

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

445 De Gruijter, J., Ter Braak, C., 1990. Model-free estimation from spatial
446 samples: A reappraisal of classical sampling theory. *Mathematical Geology* 22,
447 407–415.

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

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

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

456 Foster, S.D., Hosack, G.R., Lawrence, E., Przeslawski, R., Hedge, P., Caley,
457 M.J., Barrett, N.S., Williams, A., Li, J., Lynch, T., others, 2017. Spatially
458 balanced designs that incorporate legacy sites. *Methods in Ecology and Evolution*
459 8, 1433–1442.

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

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

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

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

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

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

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

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

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

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

482 Robertson, B., Brown, J., McDonald, T., Jaksons, P., 2013. BAS: Balanced
483 acceptance sampling of natural resources. *Biometrics* 69, 776–784.

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

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

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

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

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

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

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

500 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)
501 [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)
502 [assessment](https://www.epa.gov/national-aquatic-resource-surveys/national-results-and-regional-highlights-national-lakes-assessment).

503 Ver Hoef, J., 2002. Sampling and geostatistics for spatial data. Ecoscience 9,
504 152–161.

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

507 Ver Hoef, J.M., Temesgen, H., 2013. A comparison of the spatial linear
508 model to nearest neighbor (k-nn) methods for forestry applications. PLOS ONE
509 8, e59129.

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

513 Wang, J.-F., Stein, A., Gao, B.-B., Ge, Y., 2012. A review of spatial sampling.

514 Spatial Statistics 2, 1–14.

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