

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

Michael Dumelle<sup>\*,a</sup>, Matt Higham<sup>b</sup>, Lisa Madsen<sup>c</sup>, Anthony R. Olsen<sup>a</sup>, Jay M. Ver Hoef<sup>d</sup>

<sup>a</sup>United States Environmental Protection Agency, 200 SW 35th St, Corvallis, Oregon, 97333

<sup>b</sup>Saint Lawrence University Department of Mathematics, Computer Science, and Statistics, 23 Romoda Drive, Canton, New York, 13617

<sup>c</sup>Oregon State University Department of Statistics, 239 Weniger Hall, Corvallis, Oregon, 97331

<sup>d</sup>Marine Mammal Laboratory, Alaska Fisheries Science Center, National Oceanic and Atmospheric Administration, Seattle, Washington, 98115

## Abstract

This is the abstract.

*Text based on elsarticle sample manuscript, see <http://www.elsevier.com/author-schemas/latex-instructions#elsarticle>*

Potential Journals:

- Ecological Applications
- Methods in Ecology and Evolution
- Journal of Applied Ecology
- Environmetrics
- Environmental and Ecological Statistics

## 1. Introduction

Add sentence about bias. Add sentence about user functions.

There are two general approaches for using data to make statistical inferences about a population: design-based and model-based. When data cannot be obtained for all units in a population (population units), data on a subset of the population units is collected and called a sample. In the design-based approach, inferences about the underlying population are informed from a probabilistic process in which population units are selected to be in the sample. Alternatively, in the model-based approach, inferences are made from specific assumptions about the underlying process that generated the data. Each paradigm has a deep historical context (Sterba, 2009) and its own set of general advantages (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

---

\*Corresponding Author

Email addresses: [Dumelle.Michael@epa.gov](mailto:Dumelle.Michael@epa.gov) (Michael Dumelle), [mhigham@stlaw.edu](mailto:mhigham@stlaw.edu) (Matt Higham)  
(Manuscript submitted to An awesome journal October 19, 2021)

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. Thereafter, several comparisons between design-based and model-based for spatial data have been considered (Brus and De Gruijter, 1997; Ver Hoef, 2002, 2008). Cooper (2006) review the two approaches in an ecological context before introducing a “model-assisted” variance estimator that combines aspects from each approach. In addition to Cooper (2006), there has been substantial research and development into estimators that use both design and model-based principles (see e.g. Cicchitelli and Montanari (2012), Chan-Golston et al. (2020) for a Bayesian approach, and Sterba (2009)). More recent overviews include Brus (2020) and Wang et al. (2012).

Though comparisons between design-based and model-based approaches to spatial data have been studied, no numerical comparison has been made between design-based approaches that incorporate spatial locations and model-based approaches. In this manuscript, we compare design-based approaches that incorporate spatial locations to model-based approaches for spatial data. Though these comparisons generalize to both finite populations (e.g. point resources) and infinite populations (e.g. linear and areal resources), we focus on applications to finite populations. The rest of the manuscript is organized as follows. In Section 2, we compare sampling and estimation procedures between the design-based approach and the model-based approach for spatial data. In Section 3, we use a simulation approach to study the behavior and performance of both approaches. In Section 4, we use both approaches to analyze real data. And in Section 5, we end with a discussion and provide directions for future research.

## 2. Background

The design-based and model-based approaches incorporate randomness in fundamentally different ways. In this section, we describe the role of randomness and its effects on subsequent inferences. We then discuss specific inference methods for the design-based and model-based approaches for spatial data.

### 2.1. Comparing Design-Based and Model-Based Approaches

The design-based approach assumes the population is fixed. Randomness is incorporated via the selection of units in a sampling frame according to a sampling design. A sampling frame is the set of all units available to be sampled. A sampling design assigns a positive probability of inclusion (inclusion probability) to each unit in the sampling frame. Some examples of commonly used sampling designs include simple random sampling, stratified random sample, and cluster sampling. These sampling designs tend to select units from the sampling frame independently of other units, so we call them “Independent Random Sampling” (IRS) designs. Sampling designs incorporating the spatial locations of units in the sample frame are called spatially balanced designs. Spatially balanced designs can be obtained using the Generalized Random Tessellation Stratified

79 (GRTS) algorithm (Stevens and Olsen, 2004), which we discuss in more detail  
80 in Section 2.2. The design-based approach combines the randomness of the  
81 sampling design and the data collected via the sample to estimate parameters  
82 (e.g. means and totals) of a population. Generally, these population parameters  
83 are assumed to be fixed, unknown constants.

84 Treating the data as fixed and incorporating randomness through the sampling  
85 design yields estimators having very few other assumptions. Confidence intervals  
86 for these types of estimators are typically derived using limiting arguments  
87 that incorporate all possible randomizations of sampling units selected via the  
88 sampling design. Means and totals, for example, are asymptotically normally  
89 distributed (normal) by the Central Limit Theorem (under some assumptions).  
90 If we repeatedly sample the surface, then 95% of all 95% confidence intervals  
91 constructed from a procedure with appropriate coverage will contain the true,  
92 fixed mean. Särndal et al. (2003) and Lohr (2009) provide thorough reviews of  
93 the design-based approach.

94 The model-based approach assumes the data are a random realization of  
95 a data-generating process. Randomness is incorporated through distributional  
96 assumptions on this process. Strictly speaking, randomness need not be incor-  
97 porated through random sampling, though Diggle et al. (2010) warn against  
98 preferential sampling. Preferential sampling occurs when the process generating  
99 the data locations and the process being modeled are not independent of one  
100 another. To guard against preferential sampling, model-based approaches often  
101 still implement random sampling.

102 Instead of estimating fixed but unknown parameters (as in the design-based  
103 approach), the goal of model-based inference in the spatial context is often to  
104 predict a realized variable, or value. For example, suppose the realized mean of all  
105 population units is the value of interest. Instead of *estimating* a fixed, unknown  
106 mean, we are *predicting* the value of the mean, a random variable. Prediction  
107 intervals are then derived leveraging assumptions of the data generating process.  
108 If we repeatedly generate the response values from a fixed spatial process and  
109 obtained a sample, then 95% of all 95% prediction intervals constructed from a  
110 procedure with appropriate coverage will contain their respective realized means.  
111 Cressie (1993) and Schabenberger and Gotway (2017) provide reviews of model-  
112 based approaches for spatial data. A visual comparison of the design-based and  
113 model-based assumptions is provided in Figure 1 (Brus (2020) provides a similar  
114 figure).

## 115 2.2. Spatially Balanced Design and Analysis

116 Spatially balanced samples can be obtained using the design-based approach.  
117 Spatially balanced samples are useful because parameter estimates from these  
118 samples tend to vary less than parameter estimates from samples that are not  
119 spatially balanced (Barabesi and Franceschi, 2011; Benedetti et al., 2017; Graf-  
120 ström and Lundström, 2013; Robertson et al., 2013; Stevens and Olsen, 2004;  
121 Wang et al., 2013). The first spatially balanced sampling algorithm that saw  
122 widespread use was the Generalized Random Tessellation Stratified (GRTS)  
123 algorithm (Stevens and Olsen, 2004). To quantify the spatial balance of a

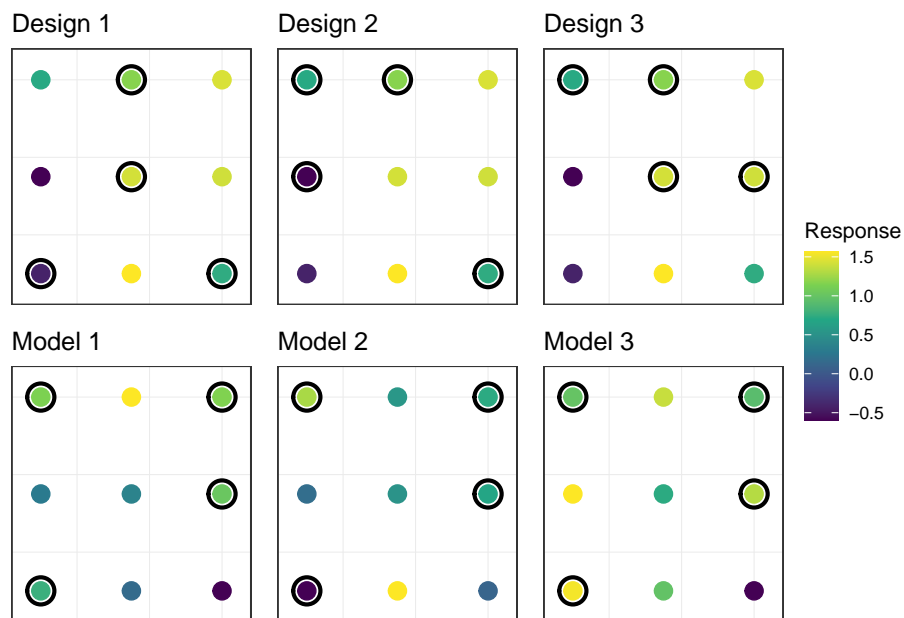


Figure 1: A comparison of sampling under the design-based and model-based frameworks. Points circled are those that are sampled. In the top row, we have one fixed population, and three random samples of  $n = 4$ . The response values at each site are fixed, but we obtain different estimates for the mean response because the randomly sampled sites vary from sample to sample. In the bottom row, we have three realizations of the same spatial process sampled at the same locations. The spatial process generating the response values has a single mean, but the realized mean is different in each of the three panels.

sample, Stevens and Olsen (2004) proposed loss metrics based on Voroni polygons. Since GRTS was developed, several other spatially balanced sampling algorithms have emerged, including the Local Pivotal Method (Grafström et al., 2012; Grafström and Matei, 2018), Spatially Correlated Poisson Sampling (Grafström, 2012), Balanced Acceptance Sampling (Robertson et al., 2013), Within-Sample-Distance Sampling (Benedetti and Piersimoni, 2017), and Halton Iterative Partitioning Sampling (Robertson et al., 2018). In this manuscript, we use Generalized Random Tessellation Stratified (GRTS) sampling because it has several attractive properties: GRTS sampling accommodates finite and infinite sampling frames; accommodates equal, unequal, and proportional (to) size inclusion probabilities; accommodates legacy (historical) sampling; accommodates a minimum distance between units in a sample; accommodates reverse hierarchically ordered replacement units in a sample (replacement units are units available to be sampled if an original unit cannot be sampled); and is available in the `spsurvey` R package Dumelle et al. (2021).

The GRTS algorithm samples from finite and infinite populations by utilizing a mapping between two-dimensional and one-dimensional space. The units in the two-dimensional sampling frame are divided into cells using a hierarchical address. This hierarchical address is then used to map the units from two-dimensional space to a one-dimensional line where each unit's line length equals its inclusion probability. A systematic sample is conducted on the line and linked back to a unit in two-dimensional space, which results in the desired sample. Stevens and Olsen (2004) and Dumelle et al. (2021) provide further details.

After selecting a GRTS sample, data are collected and used to estimate population parameters. To unbiasedly estimate population means and totals from sample data, one can use the Horvitz-Thompson estimator (Horvitz and Thompson, 1952). If  $\tau$  is a population total, the Horvitz-Thompson estimate of  $\tau$ , 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 unit in the sample and  $\pi_i$  is the inclusion probability of the  $i$ th unit in the sample. An estimate of the population mean is obtained by dividing  $\hat{\tau}_{ht}$  by the population size.

While the Horvitz-Thompson estimator is unbiased for population means and totals, it is also important to quantify the uncertainty in these estimates. Horvitz and Thompson (1952) and Sen (1953) provide variance estimators for  $\hat{\tau}_{ht}$ , but they have two drawbacks. First, these estimators rely on calculating  $\pi_{ij}$ , the probability that unit  $i$  and 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 units in the sampling frame. 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  $\pi_{ij}$  and incorporates spatial locations – for technical details see Stevens and Olsen (2003). Stevens and Olsen (2003) show the local

neighborhood variance estimator tends reduce  $\text{Var}(\hat{\tau})$  compared to variance estimators ignoring spatial locations, yielding narrower confidence intervals for  $\tau$ .

### 2.3. 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 basing inference off of a specific sampling design, we assume the data are generated by a spatial process. Ver Hoef (2008) gives details on the theory of FPBK, but some of the basic principles are summarized below. Let  $\mathbf{z} \equiv \{z(s_1), z(s_2), \dots, z(s_N)\}$  be a response vector at locations  $s_1, s_2, \dots, s_N$  that can be measured at the  $N$  population units and is represented as an  $N \times 1$  vector. Suppose we want to predict some linear function of the response variable,  $f(\mathbf{z}) = \mathbf{b}'\mathbf{z}$ , where  $\mathbf{b}'$  is a  $1 \times N$  vector of weights. For example, if we want to predict the population total across all population units, then we would use a vector of 1's for the weights.

We often only have a sample of the  $N$  population units. Denoting quantities that are part of the sampled population units with a subscript  $s$  and quantities that are part of the unsampled population units with a subscript  $u$ ,

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

where  $\mathbf{X}_s$  and  $\mathbf{X}_u$  are the design matrices for the sampled and unsampled population units, respectively;  $\beta$  is the parameter vector of fixed effects; and  $\boldsymbol{\delta}_s$  and  $\boldsymbol{\delta}_u$  are random errors for the sampled and unsampled population units, respectively. Denoting  $\boldsymbol{\delta} \equiv [\boldsymbol{\delta}_s \ \boldsymbol{\delta}_u]'$ , we assume the expectation of  $\boldsymbol{\delta}$  equals  $\mathbf{0}$ .

In addition to assuming the expectation of  $\boldsymbol{\delta}$  equals  $\mathbf{0}$ , we also assume that there is spatial correlation in  $\boldsymbol{\delta}$  that can be modeled using a covariance function. It is common to assume the covariance function is second-order stationary and isotropic (Cressie, 1993), and that the spatial covariance decreases as the separation between population units increases. Many spatial covariance functions exist, but the primary function we use throughout the simulations and applications in this manuscript is the exponential covariance function: the  $i, j$ th entry for  $\text{cov}(\boldsymbol{\delta})$  is

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

where  $\sigma_{ps}^2$  is the partial sill measuring coarse-scale (correlated) variability,  $\sigma_n^2$  is the nugget measuring fine-scale (independent) variability,  $\phi$  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  $j$ . Any spatial covariance function could be used in the place of the exponential, however, including functions that allow for non-stationarity or anisotropy (Chiles and Delfiner, 1999, pp. 80–93).

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

195 We note that we only use FPBK in this paper in order to focus more on  
 196 comparing the design-based and model-based approaches. However, k-nearest-  
 197 neighbors (Fix and Hodges, 1951; Ver Hoef and Temesgen, 2013), random  
 198 forest (Breiman, 2001), Bayesian models (Chan-Golston et al., 2020), among  
 199 others, can also be used to obtain predictions for a mean or total from spatially  
 200 correlated responses of a finite population. We choose to use FPBK because it  
 201 is faster than a Bayesian approach and random forest and because Ver Hoef and  
 202 Temesgen (2013) showed that the method outperforms k-nearest-neighbors in  
 203 many scenarios.

### 204 3. Numerical Study

205 We used a numerical simulation study to investigate performance of four  
 206 design-analysis combinations, summarized in Table 1.

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

Table 1: Types of Sampling Design and Analysis combinations considered in the simulation study. The rows give the two types of sampling designs while the columns give the two types of analyses.

207 We used a crossed design with the simulation parameters given in Table 2  
 208 for a total of 36 scenarios. All scenarios used exponential correlation with a  
 209  $\sqrt{2}/3$  for  $N = 900$  response values simulated on the unit square in either random  
 210 locations (Layout = Random) or gridded locations (Layout = Gridded). The  
 211 range was chosen so that the decayed to nearly zero at the largest distance  
 212 possible in the domain,  $\sqrt{2}$ , otherwise known as the “effective” range (for the  
 213 exponential covariance, the effective range is  $3\phi$ ). The mean for the spatial  
 214 process generating the response was set to zero.

215 For the lognormal scenarios, the response values were simulated using the  
 216 specified correlation parameters using a normal distribution and were subse-  
 217 quently exponentiated. A total variance of 2 and a mean of 0 on the normal scale  
 218 is equivalent to a total variance of 47 and a mean of 2.72 after exponentiation.  
 219 Therefore, when the model-based methods were used for lognormal response,  
 220 the correlation was mis-specified. We chose to simulate values with a lognormal  
 221 distribution so that we could test the model-based analysis approach with a  
 222 mis-specified model and so that we could test both analysis approaches on data  
 223 that exhibits a large amount of skewness.

224 There were 2000 simulation trials for each of the 36 parameter combinations.  
 225 In each trial, response values were generated from a spatial process with the

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

Table 2: Simulation parameters. Total variability for all scenarios was 2 so that the partial sill was 0, 1, or 1.8.

specified parameters, and a GRTS sample and an IRS sample were selected. For the GRTS sample, the design-based approach using the local neighborhood variance (GRTS-Design) and a model-based approach were applied (GRTS-Model). For the IRS sample, the design-based approach using the simple random sample variance (IRS-Design) and a model-based approach were applied (IRS-Model).

The GRTS algorithm and the local neighborhood variance estimator are available in the **R** package `spsurvey` (Dumelle et al., 2021). FPBK can be readily performed in **R** with the `sptotal` package (Higham et al., 2021). We use `sptotal` for both the simulation analysis and the application, estimating parameters with Restricted Maximum Likelihood (REML).

**Mike** For design-based, it’s really RMSE – how should we address this? Figure 2 shows the relative efficiency of the four approaches from Table 1 with “IRS-Design” as the baseline:

$$EFF = \frac{\text{rMSPE of approach}}{\text{rMSPE of IRS-Design}},$$

where rMSPE is the root-mean-squared-prediction error. When there is no spatial correlation (top row), the four approaches have approximately equal rMSPE, even when the assumptions of the model-based approaches are violated. So, using GRTS or using a spatial model does not result in much, if any, loss in efficiency even if the response variable is not spatially correlated. When there is high spatial correlation (bottom row), the GRTS-Model approach tends to perform best, but difference in relative efficiency between GRTS-Model and GRTS-Design is small. In the lognormal, high partial sill settings (bottom-right facet), GRTS-Design outperforms IRS-Model by a large margin, suggesting that the design decision (whether to use IRS or GRTS) is more important than the analysis decision (whether to analyze using model assumptions or not).

Unsurprisingly, Figure 2 also shows that, when the assumptions for GRTS-Model are satisfied, the approach outperforms GRTS-Design. However, even when the model that generates the data is different than the model used to fit the data, as in the lognormal response, the model-based approach still outperforms the design-based approach when there is a high amount of spatial correlation.

We also studied 95% interval coverage among the approaches. The design-based 95% confidence intervals and model-based 95% prediction intervals are constructed using the normal distribution. Justification for the design-based intervals lies in the asymptotic normality of totals via the Central Limit Theorem, and justification for the model-based intervals lies in the normality assumption



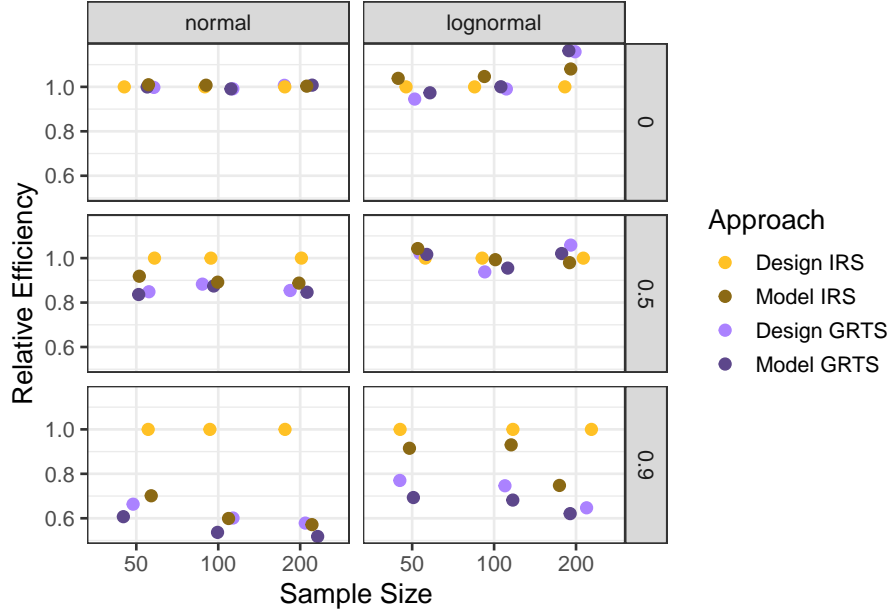


Figure 2: Relative Efficiency of the four design-analysis approaches. The plot is faceted by the type of response on the columns and the partial-sill to total-variance ratio on the rows.

of the errors. Figure 3 shows the 95% interval coverage for each of the four approaches. All four approaches have somewhat similar interval coverage in all settings, with GRTS-Design having slightly lower coverage when the response is normal.

In the normal response settings, all approaches have coverage around 95%. This is expected, as the intervals are also based on the normal distribution. In the lognormal response settings, however, all approaches have coverage below 95%. This is also expected, as the intervals are still based on the normal distribution. In the lognormal response settings, interval coverage increases both as the sample size increases and as the strength of spatial dependence increases. This suggests that the larger the sample size and the stronger the spatial dependence, the more resistant these intervals are to departures from normality of the data.

#### 4. 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 the 2012 National Lakes Assessment (NLA), which measures various aspects of lake health and quality in lakes in the contiguous United States, to obtain an interval for mean mercury concentration. Although we know the true mean mercury

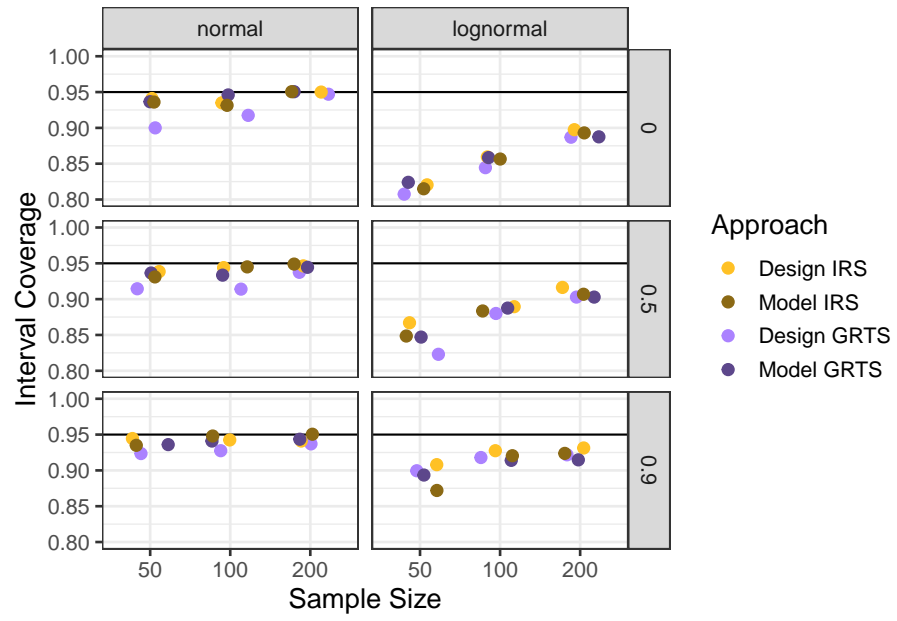


Figure 3: Coverage of the four design-analysis approaches. All confidence intervals are normal-based and have a nominal confidence level of 0.95, marked with a horizontal line. The plot is faceted by the type of response on the columns and the partial-sill to total-variance ratio on the rows.

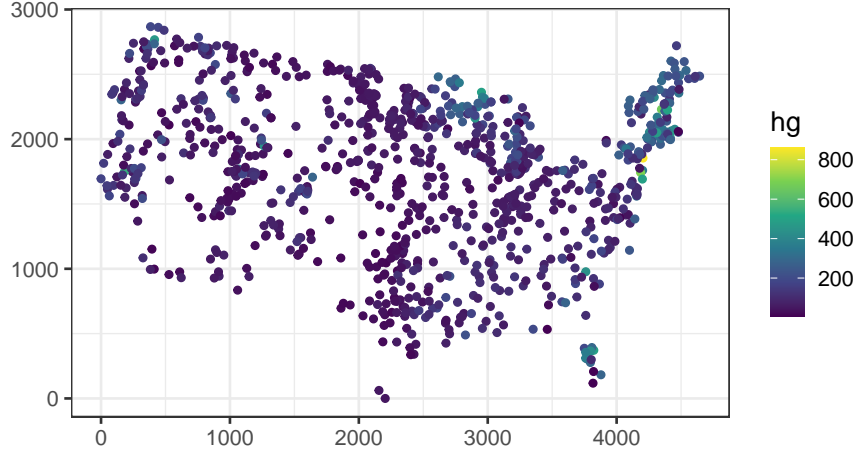


Figure 4: Population distribution of mercury concentration for 986 lakes in the contiguous United States.

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.

Figure 4 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. Mercury concentration exhibits some spatial correlation, with high mercury concentrations in lakes in the northeast and north central United States. The realized mean mercury concentration in the 986 lakes is 103.2 ng / g.

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.

Table 3 shows the application of a design-based analysis of an IRS sample, a model-based analysis of an IRS sample, a design-based analysis of a GRTS sample, and a model-based analysis of a GRTS sample. For all four analyses, the true

289 realized mean mercury concentration is within the bounds of the 95% intervals.  
 290 However, we should not generalize the results of this particular realization to  
 291 any other data set or even to other potential samples of this data set.

292 But, we do note a couple of patterns. The design-based IRS analysis shows  
 293 the largest standard error: a likely reason is that this is the only approach that  
 294 does not incorporate any spatial information regarding mercury concentration  
 295 across the contiguous United States. We also see that both approaches using  
 296 the GRTS sample have a lower standard error than the both approaches using  
 297 the IRS sample. We would expect this to be the case for most samples because  
 298 mercury concentration exhibits spatial patterning, so a spatially balanced sample  
 299 should usually yield a lower standard error.

## 300 5. Discussion

301 The design-based and model-based approaches to inference are fundamen-  
 302 tally different paradigms by which to select samples and analyze data. The  
 303 design-based approach incorporates randomness through sampling to estimate  
 304 a population parameter. The model-based approach incorporates randomness  
 305 through distributional assumptions to predict the realized value of a random  
 306 variable. Though these approaches have often been compared in the literature  
 307 both from theoretical and analytical perspectives, our contribution lies in study-  
 308 ing them in a spatial context while implementing spatially balanced sampling.  
 309 Aside from the theoretical differences described, a few analytical findings were  
 310 particularly notable: the design decision (GRTS vs IRS) seems much more im-  
 311 portant than the analysis decision (design-based vs model-based); Independent  
 312 of the analysis approach, there is no reason to prefer IRS over GRTS for spatial  
 313 data – GRTS tends to perform at least as well as IRS when there is no spatial  
 314 correlation increasingly than IRS as the strength of spatial correlation increases;  
 315 the gap in relative efficiency between GRTS-design and GRTS-model widens  
 316 as the strength of spatial correlation increases; and when the data are skewed,  
 317 interval coverage for all approaches improves both as the sample size increases  
 318 and as the strength of correlation increases.

319 There are several benefits and drawbacks of the design-based and model-based  
 320 approaches for spatial data, some of which we have not yet discussed but are  
 321 worthy of consideration in future research. The design-based approach relies on  
 322 few assumptions, while the model-based approach relies on rigid distributional  
 323 ones. The Horvitz-Thompson estimator of means and totals is unbiased and  
 324 computationally efficient, while model-based estimation of covariance parameters  
 325 can be computationally burdensome, especially for likelihood-based methods such  
 326 as REML that rely on inverting a covariance matrix. The design-based approach  
 327 also more naturally handles binary data, free from the more complicated logistic  
 328 regression formulation commonly used to handle binary data in a model-based  
 329 approach. The model-based approach, however, can quantify the relationship  
 330 between covariates (predictor variables) and the response variable, something  
 331 the design-based approach cannot do naturally. The model-based approach also  
 332 yields estimated spatial covariance parameters, which help better understand the

process of study. Model selection is also possible using model-based approaches and criteria such as likelihood ratio tests or AIC (Akaike, 1974). Model-based approaches are capable of more efficient small-area estimation than design-based approaches by leveraging distributional assumptions in areas with few observed sites. Model-based approaches can also compute site-by-site predictions at unobserved locations and use them to construct informative visualizations. The benefits and drawbacks of both approaches, alongside our theoretical and analytical comparisons, should be heavily considered when choosing among them. This is especially true from an analysis perspective, as we found that using a spatially balanced sampling algorithm benefits both design-based and model-based analyses.

## References

- Akaike, H., 1974. A new look at the statistical model identification. *IEEE transactions on automatic control* 19, 716–723.
- Barabesi, L., Franceschi, S., 2011. Sampling properties of spatial total estimators under tessellation stratified designs. *Environmetrics* 22, 271–278.
- Benedetti, R., Piersimoni, F., 2017. A spatially balanced design with probability function proportional to the within sample distance. *Biometrical Journal* 59, 1067–1084.
- Benedetti, R., Piersimoni, F., Postiglione, P., 2017. Spatially balanced sampling: A review and a reappraisal. *International Statistical Review* 85, 439–454.
- Breiman, L., 2001. Random forests. *Machine learning* 45, 5–32.
- Brus, D., De Gruijter, J., 1997. Random sampling or geostatistical modelling? Choosing between design-based and model-based sampling strategies for soil (with discussion). *Geoderma* 80, 1–44.
- Brus, D.J., 2020. Statistical approaches for spatial sample survey: Persistent misconceptions and new developments. *European Journal of Soil Science*.
- Chan-Golston, A.M., Banerjee, S., Handcock, M.S., 2020. Bayesian inference for finite populations under spatial process settings. *Environmetrics* 31, e2606.
- Chiles, J.-P., Delfiner, P., 1999. *Geostatistics: Modeling Spatial Uncertainty*. John Wiley & Sons, New York.
- Cicchitelli, G., Montanari, G.E., 2012. Model-assisted estimation of a spatial population mean. *International Statistical Review* 80, 111–126.
- Cooper, C., 2006. Sampling and variance estimation on continuous domains. *Environmetrics: The official journal of the International Environmetrics Society* 17, 539–553.
- Cressie, N., 1993. *Statistics for spatial data*. John Wiley & Sons.
- De Gruijter, J., Ter Braak, C., 1990. Model-free estimation from spatial samples: A reappraisal of classical sampling theory. *Mathematical geology* 22, 407–415.
- Diggle, P.J., Menezes, R., Su, T.-l., 2010. Geostatistical inference under preferential sampling. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 59, 191–232.

377 Dumelle, M., Olsen, A.R., Kincaid, T., Weber, M., 2021. Selecting and  
378 analyzing spatial probability samples in r using spsurvey. Manuscript Submitted  
379 for Publication.

380 Fix, E., Hodges, J.L., 1951. Discriminatory analysis, nonparametric discrimi-  
381 nation: Consistency properties. USAF School of Aviation Medicine.

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

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

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

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

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

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

395 Horvitz, D.G., Thompson, D.J., 1952. A generalization of sampling with-  
396 out replacement from a finite universe. *Journal of the American statistical*  
397 *Association* 47, 663–685.

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

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

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

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

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

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

410 Sterba, S.K., 2009. Alternative model-based and design-based frameworks  
411 for inference from samples to populations: From polarization to integration.  
412 *Multivariate behavioral research* 44, 711–740.

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

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

417 Ver Hoef, J., 2002. Sampling and geostatistics for spatial data. *Ecoscience* 9,  
418 152–161.

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

421 Ver Hoef, J.M., Temesgen, H., 2013. A comparison of the spatial linear model  
422 to nearest neighbor (k-nn) methods for forestry applications. *PloS one* 8, e59129.

423 Wang, J.-F., Jiang, C.-S., Hu, M.-G., Cao, Z.-D., Guo, Y.-S., Li, L.-F., Liu, T.-  
424 J., Meng, B., 2013. Design-based spatial sampling: Theory and implementation.  
425 Environmental modelling & software 40, 280–288.  
426 Wang, J.-F., Stein, A., Gao, B.-B., Ge, Y., 2012. A review of spatial sampling.  
427 Spatial Statistics 2, 1–14.