

1 A comparison of design-based and model-based
2 approaches for spatial data.

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

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

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

8 ^cOregon State University Department of Statistics, 239 Weniger Hall, Corvallis, Oregon,
9 97331

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

12 **Abstract**

This is the abstract.

13 *Text based on elsarticle sample manuscript, see [http://www.elsevier.com/](http://www.elsevier.com/author-schemas/latex-instructions#elsarticle)*
14 *author-schemas/latex-instructions#elsarticle*

15 Potential Journals:

- 16 • Ecological Applications
- 17 • Methods in Ecology and Evolution
- 18 • Journal of Applied Ecology
- 19 • Environmetrics
- 20 • Environmental and Ecological Statistics

21 **1. Introduction**

22 There are two general approaches for using data to make statistical inferences
23 about a population: design-based approaches and model-based approaches.
24 When data cannot be obtained for all units in a population (known as population
25 units), data on a subset of the population units is collected in a sample. In the
26 design-based approach, inferences about the underlying population are informed
27 from a probabilistic process in which population units are selected to be in the
28 sample. Alternatively, in the model-based approach, inferences are made from
29 specific assumptions about the underlying process that generated the data. Each
30 paradigm has a deep historical context (Sterba, 2009) and its own set of general
31 advantages (Hansen et al., 1983).

32 Though the design-based and model-based approaches apply to statistical
33 inference in a broad sense, we focus on comparing these approaches for spatial
34 data. We define spatial data as variables measured at specific geographic locations.
35 De Gruijter and Ter Braak (1990) give an early comparison of design-based and
36 model-based approaches for spatial data, quashing the belief that design-based

*Corresponding Author

Email addresses: Dumelle.Michael@epa.gov (Michael Dumelle), mhigham@stlaw.edu

(Manuscript received 11 June 2021)

June 11, 2021

approaches could not be used for spatially correlated data. Thereafter, several comparisons between design-based and model-based for spatial data have been considered, but they tend to compare design-based approaches that ignore spatial location in sampling to model-based approaches (Brus and De Gruijter, 1997; Ver Hoef, 2002; Ver Hoef, 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), but no numerical comparison has been made between design-based approaches that incorporate spatial sampling and model-based approaches.

The rest of this paper is organized as follows. In Section 2, we compare sampling and estimation procedures between the design-based approach and the model-based approach. In Section 3, we use simulated and real data to study the properties of parameter estimates from both approaches. And in Section 4, 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 vs. Model-Based

The design-based approach assumes the data are fixed. Randomness is incorporated in the selection of population units according to a sampling design. A sampling design assigns a positive probability of inclusion in the sample (inclusion probability) to each population unit. Some examples of commonly used sampling designs include independent random sampling (IRS), stratified random sampling, and cluster sampling. The goal is to use the sampling design and the sampled data to estimate population parameters like means and totals. These population parameters are typically assumed to be fixed but unknown.

Treating the data as fixed and incorporating randomness through the sampling design yields estimators having very few other assumptions. Confidence intervals for these types of estimators are typically derived using limiting arguments. Means and totals, for example, are asymptotically normally distributed by the Central Limit Theorem. Särndal et al. (2003) and Lohr (2009) provide thorough reviews of the design-based approach.

The model-based approach assumes the data are a random realization of a process. Randomness is often incorporated through distributional assumptions on the data-generating process. Instead of estimating fixed but unknown parameters (as in the design-based approach), the goal of model-based inference in the spatial

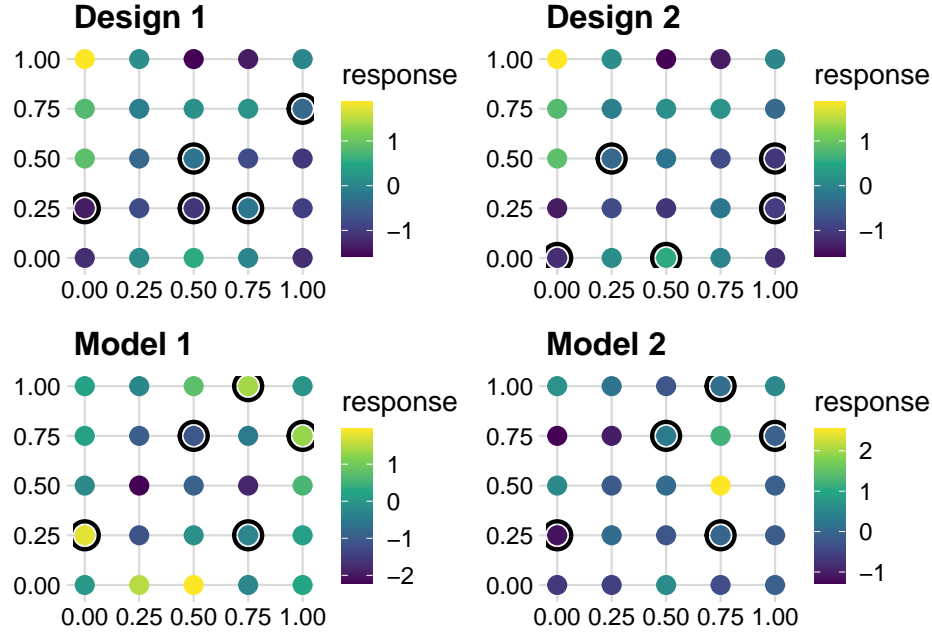


Figure 1: A comparison of sampling under the design-based and model-based frameworks. In the top row, we have one fixed population, and two random samples. In the bottom row, we have two realizations of the same spatial process sampled at the same locations.

context is often *prediction* of an unknown quantity. For example, suppose the realized mean of all population units is the quantity of interest. Instead of *estimating* a fixed unknown mean, we are *predicting* the value of the mean, a random variable. We know that if we sampled all population units, we would have an exact prediction for the mean of our one realized process, without any uncertainty. But the true mean of the spatial process that generated our realized data is still not known. When predicting the realized mean, we typically are not interested in the underlying process's true mean.

Assuming the data is a realization of a specific data-generating process yields predictors that are linked to distributional assumptions. These distributional assumptions are used to derive prediction intervals. The distributional assumptions allow the prediction intervals to be more precise. Cressie (1993) and Schabenberger and Gotway (2017) provide reviews of model-based approaches for spatial data.

2.2. Spatially Balanced Design and Analysis

Sampling designs can incorporate spatial locations to obtain samples that are spatially balanced with respect to the population (Stevens Jr and Olsen, 2004). A sample is spatially with respect to the population if the sampled population units are a miniature of the population units. A sample is a miniature of the population if the distribution of the sampled population units mirrors the density of all

99 population units. Spatial balance with respect to the population is different than
100 spatial balance with respect to geography. A sample that is spatially balanced
101 with respect to geography is spread out in some type of equidistant manner over
102 geographical space and is not meant to be miniatures of the population. When
103 we refer to spatial balance henceforth, we mean spatial balance with respect to
104 the population.

105 Spatially balanced samples are useful because they tend to yield estimates that
106 have lower variance than estimates constructed from sampling designs lacking
107 spatial balance (Barabesi and Franceschi, 2011; Benedetti et al., 2017; Grafström
108 and Lundström, 2013; Robertson et al., 2013; Stevens Jr and Olsen, 2004; Wang
109 et al., 2013). To quantify spatial balance, Stevens Jr and Olsen (2004) proposed
110 loss functions based on Voroni polygons. The first spatially balanced sampling
111 algorithm that saw widespread use was the Generalized Random Tessellation
112 Stratified (Stevens Jr and Olsen, 2004). Since GRTS was developed, several
113 other spatially balanced sampling algorithms have emerged, including the Local
114 Pivotal Method (Grafström et al., 2012; Grafström and Matei, 2018), Spatially
115 Correlated Poisson Sampling (Grafström, 2012), Balanced Acceptance Sampling
116 (Robertson et al., 2013), Within-Sample-Distance (Benedetti and Piersimoni,
117 2017), and Halton Iterative Partitioning (Robertson et al., 2018) algorithms. We
118 focus on the Generalized Random Tessellation Stratified (GRTS) algorithm to
119 select spatially balanced sampling because the algorithm has several attractive
120 properties detailed by Stevens Jr and Olsen (2004) and Dumelle et al. (2021).

121 The GRTS algorithm is used to sample from finite and infinite populations
122 and works by utilizing a mapping between two-dimensional and one-dimensional
123 space. The population units in two-dimensional space are divided into cells using
124 a hierarchical index. Population units are then mapped to a one-dimensional
125 line via the hierarchical indexing. The line length of each population unit equals
126 its inclusion probability. A systematic sample is conducted on the line and these
127 samples are linked to a population unit in two-dimensional space, which results
128 in the desired sample size. Stevens Jr and Olsen (2004) provide and Dumelle
129 et al. (2021) provide further details. The GRTS algorithm is available in the **R**
130 package **spsurvey** (Dumelle et al., 2021).

After collecting a sample, the data are used to estimate population paramet-
ters. The Horvitz-Thompson estimator (Horvitz and Thompson, 1952) and its
continuous analog (Cordy, 1993) yield unbiased estimates of population means
and totals. For example, if τ is a population total, then the Horvitz-Thompson
estimator of τ (denoted by $\hat{\tau}_{ht}$), is given by

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

131 where Z_i and π_i are the observed value and inclusion probability of the i th
132 population unit selected in the sample. Horvitz and Thompson (1952) and
133 Sen (1953) provide variance estimators for $\hat{\tau}_{ht}$, but they have two drawbacks.
134 First, they rely on calculating π_{ij} , the probability that population unit i and
135 population unit j are included in the sample, which can be very difficult to

136 calculate. Second, they ignore the spatial locations of the population units.
 137 To address these drawbacks, Stevens Jr and Olsen (2003) proposed a local
 138 neighborhood variance estimator. The local neighborhood variance estimator
 139 does not rely on π_{ij} , and it incorporates spatial locations by assigning higher
 140 weights to nearby observations. Stevens Jr and Olsen (2003) show this variance
 141 estimator tends to reduce the variability associated with estimating τ . This
 142 yields confidence intervals for τ that are narrower than confidence intervals
 143 constructed from variance estimators ignoring spatial locations.

144 2.3. Finite Population Block Kriging

145 Finite Population Block Kriging (FPBK) is an alternative to sampling-based
 146 methods (Ver Hoef, 2008). FPBK expands the geostatistical kriging framework
 147 to the finite population setting. Instead of basing inference off of a specific
 148 sampling design, we assume the data are generated by a spatial process with
 149 parameters that can be estimated using the framework of a model.

150 Ver Hoef (2008) gives details on the theory of FPBK, but some of the basic
 151 principles are summarized below. For a response variable \mathbf{z} that can be measured
 152 on a finite number of N sites, we want to predict some linear function of the
 153 response variable, $\tau(\mathbf{z}) = \mathbf{b}'\mathbf{z}$, where \mathbf{b} is a vector of weights. For example, if we
 154 want to predict the population total across all sites, then we would use a vector
 155 of 1's for the weights.

156 Typically, however, we only have a sample of the N sites. Denoting quantities
 157 that are part of the sampled sites with a subscript s and quantities that are part
 158 of the unsampled sites with a subscript u ,

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

159 where \mathbf{X}_s and \mathbf{X}_u are the design matrices for the sampled and unsampled sites,
 160 respectively, and $\boldsymbol{\delta}_s$ and $\boldsymbol{\delta}_u$ are random errors for the sampled and unsampled
 161 sites. Denoting $\boldsymbol{\delta} \equiv [\boldsymbol{\delta}_s \ \boldsymbol{\delta}_u]'$, we assume that $E(\boldsymbol{\delta}) = \mathbf{0}$.

We also typically assume that there is spatial correlation in $\boldsymbol{\delta}$, which can be
 modeled using a covariance function. Many common choices for this function
 assume that spatial covariance decreases with increasing Euclidean distance
 between sites. The primary function used throughout the simulations and
 applications of this manuscript is the Exponential covariance function: the i, j^{th}
 entry for $\text{var}(\boldsymbol{\delta})$ is

$$\text{cov}(\delta_i, \delta_j) = \theta_1 \exp(-3h_{i,j}/\theta_2) + \theta_3 \mathbb{1}\{\mathbf{h}_{i,j} = 0\}, \quad (3)$$

162 where $h_{i,j}$ is the distance between sites i and j , and $\boldsymbol{\theta}$ is a vector of spatial
 163 covariance parameters of the partial sill θ_1 , the range θ_2 , and the nugget θ_3 .
 164 However, any spatial covariance function could be used in the place of the
 165 Exponential, including functions that allow for anisotropy [pg. 80 - 93](Chiles
 166 and Delfiner, 1999).

167 With the above model formulation, the Best Linear Unbiased Predictor
 168 (BLUP) for $\tau(\mathbf{b}'\mathbf{z})$ 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. Neither require a particular distribution for \mathbf{z} .

We note that we only use FPBK in this paper in order to focus more on comparing the design-based and model-based approaches. However, k-nearest-neighbors (Fix and Hodges, 1951; Ver Hoef and Temesgen, 2013), random forest (Breiman, 2001), Bayesian models (Chan-Golston et al., 2020), among others, can also be used to obtain predictions for a mean or total from spatially correlated responses in a finite population setting.

3. Numerical Study

Sample Simulation

For the following simulation results, we simulated 1040 different gridded populations, each of size 900 with sample size 150. For the model-based approach (FPBK), sites were selected via Independent Random Sample. For GRTS, the local mean variance was used.

The response was normally distributed with an exponential covariance function with partial sill of 0.9, effective range of $\sqrt{2}$, and a nugget of 0.1. For model-based, we assumed the correct form of the covariance function (Exponential), but estimated the spatial parameters with REML.

Approach	Bias	RMSE	MedAE	Coverage	PClose	MedIL
Design	0.0003	0.0353	0.0251	0.9461	0.4889	0.1362
Model	-0.0001	0.0362	0.0253	0.9480	0.5111	0.1430

Table 1: Approach, mean bias (Bias), root-mean-squared error (RMSE), median absolute error (MedAE), 95 percent interval coverage (Coverage), proportion of times the approach was closer to the true value (PClose), and median interval length (MedIL)

Base Simulations

- both good: correctly specified model with high correlation
- break model: highly non-normal errors with small sample size
- break design: small area estimation

Simulation Discussion Questions

- model-based: how should sample be drawn? should locations be fixed?
- change n or sampling fraction?

Other Base Settings?

- both good?: misspecified covariance model with high correlation
- break both? non-gaussian areas with smaller sample size

3.1. Software

FPBK can be readily performed in R with the `sptotal` package (Higham et al., 2020). We use `sptotal` for both the simulation analysis and the application, estimating parameters with Restricted Maximum Likelihood (REML).

201 3.2. *Applied Example*

202 Potential Data Sets:

- 203 • National Lakes Assessment
- 204 • Moose in Alaska
- 205 • Temperature Data from NOAA

206 4. Discussion

207 References

- 208 Barabesi, L., Franceschi, S., 2011. Sampling properties of spatial total estimators
209 under tessellation stratified designs. *Environmetrics* 22, 271–278.
- 210 Benedetti, R., Piersimoni, F., 2017. A spatially balanced design with probability
211 function proportional to the within sample distance. *Biometrical Journal* 59,
212 1067–1084.
- 213 Benedetti, R., Piersimoni, F., Postiglione, P., 2017. Spatially balanced sampling:
214 A review and a reappraisal. *International Statistical Review* 85, 439–454.
- 215 Breiman, L., 2001. Random forests. *Machine learning* 45, 5–32.
- 216 Brus, D., De Gruijter, J., 1997. Random sampling or geostatistical modelling?
217 Choosing between design-based and model-based sampling strategies for soil
218 (with discussion). *Geoderma* 80, 1–44.
- 219 Brus, D.J., 2020. Statistical approaches for spatial sample survey: Persistent
220 misconceptions and new developments. *European Journal of Soil Science*.
- 221 Chan-Golston, A.M., Banerjee, S., Handcock, M.S., 2020. Bayesian inference for
222 finite populations under spatial process settings. *Environmetrics* 31, e2606.
- 223 Chiles, J.-P., Delfiner, P., 1999. *Geostatistics: Modeling Spatial Uncertainty*.
224 John Wiley & Sons, New York.
- 225 Cicchitelli, G., Montanari, G.E., 2012. Model-assisted estimation of a spatial
226 population mean. *International Statistical Review* 80, 111–126.
- 227 Cooper, C., 2006. Sampling and variance estimation on continuous domains.
228 *Environmetrics: The official journal of the International Environmetrics*
229 *Society* 17, 539–553.
- 230 Cordy, C.B., 1993. An extension of the horvitz—thompson theorem to point
231 sampling from a continuous universe. *Statistics & Probability Letters* 18,
232 353–362.
- 233 Cressie, N., 1993. *Statistics for spatial data*. John Wiley & Sons.
- 234 De Gruijter, J., Ter Braak, C., 1990. Model-free estimation from spatial samples:
235 A reappraisal of classical sampling theory. *Mathematical geology* 22, 407–415.
- 236 Dumelle, M., Olsen, A.R., Kincaid, T., Weber, M., 2021. Selecting and analyzing
237 spatial probability samples in r using spsurvey. *Manuscript Submitted for*
238 *Publication*.
- 239 Fix, E., Hodges, J.L., 1951. Discriminatory analysis, nonparametric discrimina-
240 tion: Consistency properties. *USAF School of Aviation Medicine*.
- 241 Grafström, A., 2012. Spatially correlated poisson sampling. *Journal of Statistical*
242 *Planning and Inference* 142, 139–147.

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

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

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

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

252 Higham, M., Ver Hoef, J., Bryce, F., 2020. Sptotal: Predicting totals and
253 weighted sums from spatial data.

254 Horvitz, D.G., Thompson, D.J., 1952. A generalization of sampling without
255 replacement from a finite universe. *Journal of the American statistical*
256 *Association* 47, 663–685.

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

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

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

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

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

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

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

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

274 Stevens Jr, D.L., Olsen, A.R., 2004. Spatially balanced sampling of natural
275 resources. *Journal of the american Statistical association* 99, 262–278.

276 Ver Hoef, J., 2002. Sampling and geostatistics for spatial data. *Ecoscience* 9,
277 152–161.

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

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

283 Wang, J.-F., Jiang, C.-S., Hu, M.-G., Cao, Z.-D., Guo, Y.-S., Li, L.-F., Liu, T.-J.,
284 Meng, B., 2013. Design-based spatial sampling: Theory and implementation.
285 *Environmental modelling & software* 40, 280–288.

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