

1  
2

3  
4

5

6  
78  
9

10  
11

## 12

13

14

15

- 16  
17  
18  
19  
20

## 21

22  
23  
24  
25  
26  
27  
28  
29  
30  
31

32  
33  
34  
35  
36

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 locations to model-based approaches (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), but no numerical comparison has been made between design-based approaches that incorporate spatial locations 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 the behavior of 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 data-generating process. Randomness is often incorporated through distributional assumptions on this process. Instead of estimating fixed but unknown parameters (as in the design-based approach), the goal of model-based inference

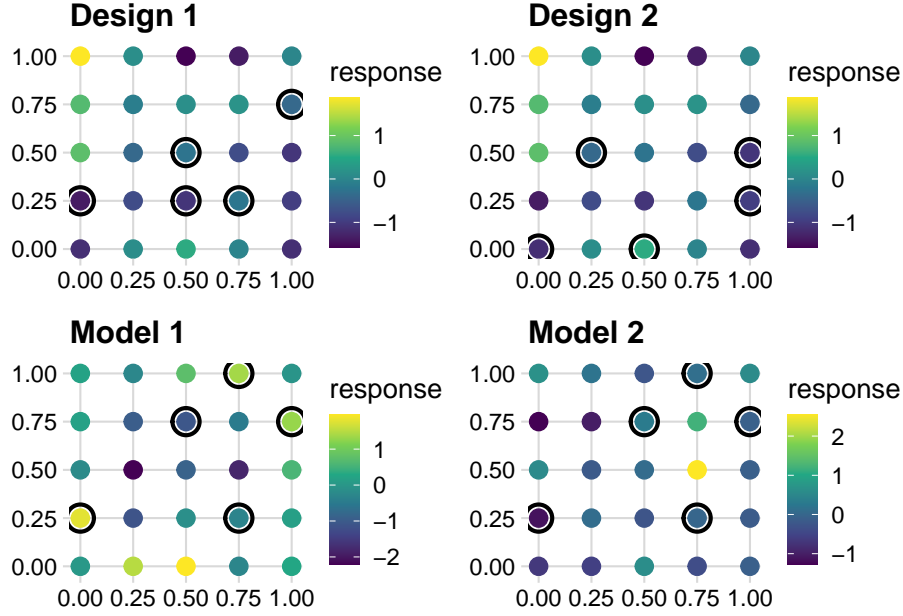


Figure 1: A comparison of sampling under the design-based and model-based frameworks. In the top row, we have one fixed population, and two random samples. In the bottom row, we have two realizations of the same spatial process sampled at the same locations.

in the spatial 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.

Figure 1 shows

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

Spatially balanced samples are useful because they tend to yield estimates that have lower variance than estimates constructed from sampling designs lacking spatial balance (Barabesi and Franceschi, 2011; Benedetti et al., 2017; Grafström and Lundström, 2013; Robertson et al., 2013; Stevens Jr and Olsen, 2004; Wang et al., 2013). To quantify spatial balance, Stevens Jr and Olsen (2004) proposed loss functions based on Voroni polygons. The first spatially balanced sampling algorithm that saw widespread use was the Generalized Random Tessellation Stratified (Stevens Jr and Olsen, 2004). 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 (Benedetti and Piersimoni, 2017), and Halton Iterative Partitioning (Robertson et al., 2018) algorithms. We focus on the Generalized Random Tessellation Stratified (GRTS) algorithm to select spatially balanced sampling because the algorithm has several attractive properties detailed by Stevens Jr and Olsen (2004) and Dumelle et al. (2021).

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

After collecting a sample, the data are used to estimate population parameters. 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  $\tau$  is a population total, then the Horvitz-Thompson estimator 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  $Z_i$  and  $\pi_i$  are the observed value and inclusion probability of the  $i$ th population unit selected in the sample. Horvitz and Thompson (1952) and Sen (1953) provide variance estimators for  $\hat{\tau}_{ht}$ , but they have two drawbacks. First, they rely on calculating  $\pi_{ij}$ , the probability that population unit  $i$  and

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

### 145 2.3. Finite Population Block Kriging

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

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

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

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

160 where  $\mathbf{X}_s$  and  $\mathbf{X}_u$  are the design matrices for the sampled and unsampled sites,  
 161 respectively, and  $\boldsymbol{\delta}_s$  and  $\boldsymbol{\delta}_u$  are random errors for the sampled and unsampled  
 162 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

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

163 where  $h_{i,j}$  is the distance between sites  $i$  and  $j$ , and  $\boldsymbol{\theta}$  is a vector of spatial  
 164 covariance parameters of the partial sill  $\theta_1$ , the range  $\theta_2$ , and the nugget  $\theta_3$ .  
 165 However, any spatial covariance function could be used in the place of the  
 166 Exponential, including functions that allow for anisotropy [pg. 80 - 93](Chiles  
 167 and Delfiner, 1999).

168 With the above model formulation, the Best Linear Unbiased Predictor  
 169 (BLUP) for  $\tau(\mathbf{b}'\mathbf{z})$  and its prediction variance can be computed. While details  
 170 of the derivation are in (Ver Hoef, 2008), we note here that the predictor and its  
 171 variance are both moment-based. Neither require a particular distribution for  $\mathbf{z}$ .

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

### 178 3. Numerical Study

#### 179 Sample Simulation

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

184 The response was normally distributed with an exponential covariance func-  
 185 tion with partial sill of 0.9, effective range of  $\sqrt{2}$ , and a nugget of 0.1. For  
 186 model-based, we assumed the correct form of the covariance function (Exponen-  
 187 tial), 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)

#### 188 Base Simulations

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

#### 192 Simulation Discussion Questions

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

#### 195 Other Base Settings?

- 196 • both good?: misspecified covariance model with high correlation
- 197 • 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).

### 3.2. Applied Example

Potential Data Sets:

- National Lakes Assessment
- Moose in Alaska
- Temperature Data from NOAA

## 4. Discussion

## References

- 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.
- Cordy, C.B., 1993. An extension of the horvitz—thompson theorem to point sampling from a continuous universe. *Statistics & Probability Letters* 18, 353–362.
- 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.

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

279 Ver Hoef, J., 2002. Sampling and geostatistics for spatial data. *Ecoscience* 9,  
280 152–161.

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

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



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