

6 Spatially balanced survey designs for natural resources

Anthony R. Olsen, Thomas M. Kincaid, and Quinn Payton

Introduction

A common objective for a monitoring program is to characterize an environmental resource based on inference from the sites selected to be monitored to the entire target population. The scale of monitoring ranges from local studies, to regional monitoring, to nationwide monitoring programs. Rarely can these monitoring efforts monitor at all locations, or sites, within the study region. Consequently, a major consideration is how to select representative sites from which it is possible to make inferences to the entire study region.

In addition to the different spatial scales of interest across monitoring efforts, these studies may focus on different environmental resources. A state or province may be interested in all small lakes (e.g. < 10 ha) with the objective of classifying the lakes as meeting designated uses (i.e. having values of designated water quality attributes that do not exceed a specified threshold), partially meeting designated uses or not meeting designated uses. In this case, the elements of the environmental resource are individual lakes and site selection is based on selecting a subset of lakes from the target population. Alternatively, a state may be interested in all perennial streams and rivers within the state to determine the total stream length that meets a nutrient criterion. In this case, the elements of the environmental resource are all possible locations on the stream and river network within the state and site selection is based on selecting sites on the stream network to be monitored. The stream network is the target population and is viewed as a continuous linear network. Finally, a state may be interested in monitoring a single large estuary within the state (e.g. Puget Sound in Washington State, USA) to determine the proportion of the estuarine area that has sediment contamination exceeding a criteria. In this case, the environmental resource, i.e. the target population, is an area and sites are selected from all possible locations within the estuarine surface area. Similar examples can be given for terrestrial environmental resources.

Some examples of US national-scale environmental resource monitoring programs are:

- National Aquatic Monitoring Surveys (NARS) conducted by the US Environmental Protection Agency and states to monitor the condition of lakes, reservoirs, streams,

Design and Analysis of Long-term Ecological Monitoring Studies, ed. R.A. Gitzen, J.J. Millsaugh, A.B. Cooper and D.S. Licht. Published by Cambridge University Press. © Cambridge University Press 2012.

rivers, wetlands, estuaries and near shore coastal waters (<http://water.epa.gov/type/watersheds/monitoring/nationalsurveys.cfm>).

- Forest Inventory and Analysis (FIA) program of the US Forest Service focusing on the Nation's forests (<http://fia.fs.fed.us/>) (discussed in Chapter 2).
- National Wetland Status and Trends program of the US Fish and Wildlife Service monitors wetland acreage changes in the conterminous United States (<http://www.fws.gov/wetlands/StatusAndTrends/index.html>).
- Vital Signs Monitoring program of the National Park Service to track the overall condition of natural resources in parks (<http://science.nature.nps.gov/im/>) (Chapter 22).

In all of these cases and in all other environmental resource monitoring programs, an important feature of their respective environmental resources is that they are distributed over geographic space. When probability sampling designs explicitly use spatial location in the selection of the sites, we define the resulting survey design as “spatially explicit”, or simply “spatial”, survey designs. These designs must be able to address the large variation in scale and the three fundamental types of target populations – those that consist of points, lines, or polygons.

This chapter first discusses (i) the concept of representative samples in this spatial context, building on the broad overview of spatial sampling provided in Chapter 5. Our chapter then (ii) describes a specific type of spatially balanced survey design, generalized random tessellation stratified (GRTS) survey design (Stevens and Olsen 2004); (iii) introduces a definition for spatial balance and proposes spatial balance metrics; and (iv) compares the spatial balance properties of simple random sampling to a GRTS survey design, using these spatial balance metrics.

Representative sample and representative sampling process

The scientific literature and specifically the environmental monitoring literature include numerous uses of the expression “representative sample”, or expressions similar to “the sample is representative of the population”. In many cases the meaning of these expressions is not explained and it is left up to the reader to interpret what the authors meant. Collectively, scientists have a general sense of what is meant when they read “representative sample”, although on closer examination the term is found to have multiple interpretations. In a series of three landmark papers in 1979, Kruskal and Mosteller (1979a, 1979b, and 1979c) reviewed the non-scientific writing and the scientific (excluding statistics) and statistical literature to classify and illustrate the various meanings of “representative sample” and “representative sampling”. They classified the meanings for “representative sampling” in the statistical literature as follows:

- General acclaim for data.
- Absence of selective forces.
- **Miniature of the population.**
- Typical or ideal case(s).

- **Coverage of the population.**
- Vague term, to be made precise.
- **Representative sampling as a specific sampling method.**
- Representative sampling as permitting good estimation.
- Representative sampling as good enough for a particular purpose.

The first six meanings also occur in non-scientific and scientific excluding statistics literature. We focus on the meanings (in bold above) that are particularly relevant to monitoring design.

What is a “miniature of the population?” Conceptually, it is a small portion of the population that arises from a perfect mixing of the population. An example is a sample of blood taken during a physical examination to ascertain particular chemical concentrations in a person’s blood. If the population had 2 categorical variables, such as sex, 2 age groups, and 2 income groups, then a miniature of the population would have the same percent of individuals in the $2 \times 2 \times 2 = 8$ categories as the entire population.

A representative sample as a miniature of the population is appealing in monitoring because the characteristics of the sample then would apply to the entire resource being monitored. Constructing a sample that is a miniature of an environmental resource is extremely difficult unless the number of factors involved is very small. For example, in monitoring lakes, factors that may be considered could be 5 lake area categories, 2 origin categories (natural versus human-made) and 3 elevation categories. In this example, all possible combinations of these factors generate 30 categories, assuming all combinations are present in the target population. A miniature of the population would require at least a sample size of 30 in order to just have a single site in each combination. A single site in each combination is unlikely to be a miniature of the population since the number of lakes in each combination is unlikely to be the same. Hence, the number of sampled sites would need to be proportional to the number of lakes in each combination which may require a very large sample size if the number of lakes in the categories are very different. Even if the sample size is sufficient to select the number of lakes required for achieving proportional allocation of lakes in the sample, a procedure for selecting lakes within each category is required.

One approach would be to use professional judgment to select representative sites within each combination. Generally samples selected in this way fall short of being miniatures due to bias in professional judgment occurring (known or unknown), i.e. selection bias (Chapters 2, 5, 11). In the context of representative sampling as a miniature of the population Kruskal and Mosteller (1979c: 251) quote the following definition of representative sample by Stephan and McCarthy (1958: 31–32): “A *representative sample* is a sample which for a specified set of variables, resembles the population. . . . [in that] certain specified analyses. . . . (computation of means, standard deviations, etc. . . .) will yield results. . . . within acceptable limits set about the corresponding population values, except that. . . . [rarely] the results will fall outside the limits . . . the mere statement or claim that a sample is representative of a population tells us nothing.” A representative

sample that is a miniature of the population is useful when inferences to the entire population are the focus. Thus the issue is: how do we know that the sites are a miniature of the population?

Coverage of the population is the concept that a representative sample is best achieved by getting the sample from as many parts, or partitions, of the population as possible. In that sense it is similar to the concept of a miniature of the population. The difference is that coverage of the population does not imply that the sample has the same relative frequencies in the sample as the population as does the “miniature of the population” concept of representative sampling. This type of representative sampling may be useful in the context of building models where the partitions ensure that the range of each variable being considered in the model is in the sample. When the monitoring objective is inference to the entire population, a representative sample that only provides coverage of the population is not sufficient.

A third meaning is representative sampling as a specific sampling method. Note that this meaning focuses on a process to generate a sample – not to a specific sample. Fundamental to this meaning is that each element of the population has a chance of appearing in any realization of the representative sampling process. The simplest representative sampling process is a simple random sample where each element of the population has an equal chance of being selected. For example, in a study of all lakes (N of them) in the State of Oregon, USA, each lake would have a chance of being selected as one of the lakes in a sample of size n . In this case each lake has a probability of n/N of being selected in the sample. This is a simple example of what is called a probability survey design (see Chapters 2, 5). Kruskal and Mosteller (1979c: 257) recognize that while simple random samples are intuitively representative in a long-run average sense, specific realizations can be “wildly unrepresentative”. Note that the likelihood of “wild realizations” decreases as the sample size becomes larger. For example, a single realization of a simple random sample of lakes in Oregon may lead to all lakes being located in western Oregon and include only lakes < 10 ha. It is unlikely that such a sample would be accepted as a representative sample of lakes in Oregon. That is, a representative sampling process does not necessarily produce a representative sample.

The appeal of a representative sampling process is that statistical inferences are then available that allows quantification of the uncertainty associated with any estimates (Chapters 2, 5). Snedecor and Cochran (1967) define probability sampling (i.e. a representative sampling process) as having the properties that (i) every unit of the population has a known probability (> 0) of being included in the sample; (ii) the sample is selected by some method of random selection, or systematic selection or systematic selection with a random start, consistent with probabilities; and (iii) the probabilities of selection (weighting factors) are used in making inferences from the sample to the population in question. Note that representative sampling processes may purposely select samples that are not miniatures of the population (e.g. by using stratification) and then account for this when statistical estimation is completed. In this case of stratification, the representative sampling process applied in each stratum is intended to result in a representative sample that is a miniature of the population within the stratum.

Obtaining a representative sample in environmental monitoring

In monitoring designs, what can be done when implementing a representative sampling process to improve the chance that each realization of the process results in a representative sample, i.e. is a miniature of the population? One approach is to stratify by characteristics of the population, then sample each stratum proportional to its occurrence in the population and within each stratum select sites with equal probability. For example, lakes may be stratified by four classes based on lake surface area, four categories of ecoregions, three elevation categories, and as natural versus man-made lakes. The stratification guarantees that the sample will include samples from each of the strata. The proportional sampling ensures that each stratum is sampled in proportion to its occurrence in the population. The equal probability of selection reduces the chance of selection bias within a stratum. However, stratification requires information upon which the strata can be defined for the entire population. Such information may not be readily available, or the number of strata desired may make it impossible to proportionally allocate the sample to the strata given the sample size available. For example, natural versus man-made lake information is typically not available. For the example, a total of 96 strata are required which will make it impossible to proportionally allocate a sample size of 100 to the strata.

Another approach is to construct strata that partition the population spatially so that the sample is a miniature of the population spatially. Constructing explicit spatial strata and sampling them proportional to the number of lakes in each spatial stratum also leads to the same problem when the number of spatial strata is large. An alternative way to use spatial strata underlies the spatial survey design approach described in this chapter.

Our view is that environmental monitoring designs should incorporate the concept of a representative sample being a miniature of the population and that the representative sample should be the result of a representative sampling process (Box 6.1). In addition, since environmental resources are intrinsically spatial in nature, monitoring designs should explicitly use spatial location in the selection of the sites. Conceptually, we postulate that spatial distribution is a useful surrogate to using combinations of characteristics of the population to get a representative sample that is a miniature of the population. That is, *representative samples with the same spatial distribution as the population are more likely to be miniatures of the population than samples that do not have same spatial distribution as the population*. Space is not always sufficient on its own. For example, it is possible that a lake sample that has the same spatial distribution as the population could consist of only small lakes. Our contention is that spatial survey designs are more likely to minimize these occurrences than not. Spatial survey designs incorporate these principles.

To illustrate, assume we have a sample of lakes in Oregon that has the same spatial distribution as all lakes in Oregon. Intuitively, the perception is that if the sample has the same spatial distribution as the population of lakes then the sample is more likely a miniature of the population. Our perception that the sample is a miniature of the lake population is likely increased if we know that the sample is the result of a representative sampling process. The representative sampling process reduces the

Box 6.1 Take-home messages for program managers

The “GRTS” spatially explicit survey designs for natural resources were created to address several issues that are commonly faced when a monitoring program based on probability surveys is designed. First, natural resources occur as discrete objects (e.g. whole lakes), linear networks (e.g. streams), or collections of areas (e.g. areas within a collection of habitat patches), represented in Geographic Information Systems (GIS) as collections of points, linear networks, or polygons. The GRTS designs have options for sampling any of these three types. Second, the GRTS spatial survey design process is a representative sampling process that is constrained so that each potential sample generated is a representative sample, a miniature of the natural resource in space. That is, each sample has the same spatial distribution as the natural resource population. While this has significant statistical advantages compared to alternatives such as simple random or systematic sampling, it also has the advantage that when the sample is displayed graphically viewers are likely to accept the sample as being representative – and therefore more likely to accept the value of the estimates of status or trends resulting from the sample. The availability of a software program to create GRTS spatial survey designs makes these designs easily accessible to those designing monitoring programs. Hundreds of such designs have been created at local to national scales and for aquatic and terrestrial natural resources (e.g. Chapters 10, 17). Our recommendation is that GRTS designs should be the default choice whenever probability sampling is used in monitoring programs.

Implementing a GRTS spatial survey design, regardless of the design’s complexity, requires that statistical analyses using data from the design incorporate the properties of the design, e.g. stratification or unequal probability of selection. That is, the statistical analysis must match the spatial survey design. This may be a change in culture for an organization using a spatially explicit survey design for the first time. Selecting the design depends on the monitoring program having a set of clearly defined, quantitative objectives (Chapters 2, 3, 18). Developing these objectives is typically the most difficult task in designing the monitoring program. It requires that managers and those developing the monitoring program discuss and agree on exactly what the program will produce.

likelihood of selection bias and increases our perception that characteristics of the lakes that are important to the study will be proportionally represented in the sample of lakes.

The concept of incorporating spatial regularity in sampling environmental populations is well established. Systematic sampling using a regular grid (or variations thereof) has been studied and used extensively (Bickford *et al.* 1963, Olea 1984, Gilbert 1987; Chapter 5). These are common in studies of areal resources such as forests or contaminated areas. Stevens (1997) developed generalizations of grid-based systematic designs, i.e. random tessellation stratified (RTS) designs, and applied them to areal resources. Subsequently, Stevens and Olsen (2004) generalized the RTS designs, making them applicable to point and line environmental resources as well as areal resources. These

generalized random tessellation stratified (GRTS) designs are spatial survey designs that explicitly incorporate the concept of spatial balance as a means of increasing the likelihood that each realization of the representative sampling process results in a miniature of the population. A GRTS spatial survey design is a representative sampling process that results in representative samples that are more likely to be miniatures of the population than other common spatial survey designs.

Generalized random tessellation stratified (GRTS) survey designs

Stevens and Olsen (2004) present the theory behind GRTS designs and give an example of its application. A software implementation of the GRTS algorithm is available in the “spsurvey” package [function *grts()*] developed for the R statistical software (R Core Development Team, 2010). The “spsurvey” package is available on the R CRAN website (<http://www.r-project.org/>) and the US Environmental Protection Agency’s Aquatic Resource Monitoring (ARM) website (<http://www.epa.gov/nheerl/arm/>). The latter site also includes information on designing monitoring programs for aquatic resources. Our objective in this section is to describe the GRTS algorithm in accessible terms to users. Use of the “spsurvey” *grts()* function is demonstrated in Appendix 6.1. Readers should also consider additional implementation issues discussed in Box 6.2.

Basic GRTS algorithm: GRTS equal probability spatial survey design

To describe the GRTS spatially balanced algorithm, we use a simple illustrative example of selecting a sample of size 8 from a population of 1957 lakes in Oregon. In this example, the primary monitoring objective is to estimate the current status of the lakes in terms of a benthic macroinvertebrate index of biotic integrity (IBI; Stoddard *et al.* 2008) and report on the number or proportion of lakes having specific ranges of the IBI score. Conceptually, each lake could have a benthic macroinvertebrate sample collected for the lake using an appropriate field design, taxonomic identification could be completed and a value for the index computed for the lake. In this case, lakes are considered a point environmental resource because the objective is to report on the number of lakes in terms of their biological integrity where biological integrity is based on a single lake-wide measure of integrity. Assume the target population is all lakes that are greater than 1 ha, more than 1 m deep, and with at least 10% of their surface as open water. The GRTS algorithm requires a list of all lakes in the target population along with their location as *x*-,*y*-coordinates, usually in an equal area projection (e.g. Albers) rather than geographic coordinates. We use the lakes within Oregon which are greater than 1 ha in the National Hydrography Dataset (<http://nationalmap.gov/>) as that list, i.e. the sample frame from which the lakes will be selected.

The core concept of the GRTS algorithm is the creation of spatial strata by constructing a grid that satisfies specific requirements. The following describes what the algorithm and software that implement it does. A map showing the sample frame with a square grid overlaid is shown in Fig. 6.1a. The grid is constructed such that (i) the number

Box 6.2 Common challenges: implementing GRTS

The main focus of this chapter is providing the conceptual basis of spatial survey designs and their implementation using the GRTS process applied to natural resources modeled as points, linear networks or polygons. The implementation using the “spsurvey” package in R requires the preparation of sample frame as a GIS layer which is imported to R as an ESRI shapefile (Appendix 6.1). This requires users to have access to ArcGIS, which is less of an issue than a few years ago. However, the creation of a single GIS layer that has the attributes required for the spatial survey design and identifies the objects in the layer that should be included in the sample frame typically requires more effort than users expect.

The first major decision is whether to model the natural resource as a point, linear network or polygons based on the monitoring objectives. For example, the user may have a polygon lake GIS layer but the monitoring objectives may require that lakes be modeled as points. In this case the polygon layer must be converted to a point layer. The coordinate system used by the GIS layer also matters. In most cases, a GIS layer in geographic coordinates (latitude, longitude) is not appropriate for use in spatial survey design. The reason is that the distance represented by a degree of longitude is not the same distance represented by a degree of latitude. Typically, an area preserving projection should be used, e.g. an Albers or UTM projection).

A monitoring program based on a complex spatial survey design also requires using an appropriate statistical analysis of the monitoring data. Not doing so can result in incorrect estimates. Even when the spatial survey design is a non-stratified, equal probability design, the statistical analysis can be improved by using more complex variance estimation procedures. Since spatially balanced designs reduce the probability of unusual samples, a variance estimator can be defined that is unbiased and that produces smaller estimates than the usual simple random sample variance estimator. Stevens and Olsen (2003) defined this local neighborhood variance estimator and showed that it performs better than alternatives. Unfortunately, the estimator cannot be computed without the use of the R software (or other software that computes generalized inverses for matrices). For spatial survey design data, the “spsurvey” package includes functions that will calculate means, totals, percentiles and cumulative distribution functions as well as their local neighborhood variances.

of rows and columns must be a power of 2, i.e. 2^m where $m > 0$ and in this case is 2, and (ii) it has an extra row and column of empty grid cells at the top and right. The choice of m is discussed later, although it depends on the sample size desired and the spatial distribution of lakes. The user does not need to specify m as the GRTS algorithm computes it. The extra row and column are used to ensure that any two lakes may occur in different cells. This is accomplished by randomly shifting the sample frame by selecting a single random x and random y from uniform distributions over their ranges defined by the boundaries of the lower left cell and moving all lake (i.e. the sample frame) by adding the random shifts to their coordinates. Figure 6.1b shows the resulting shift. The

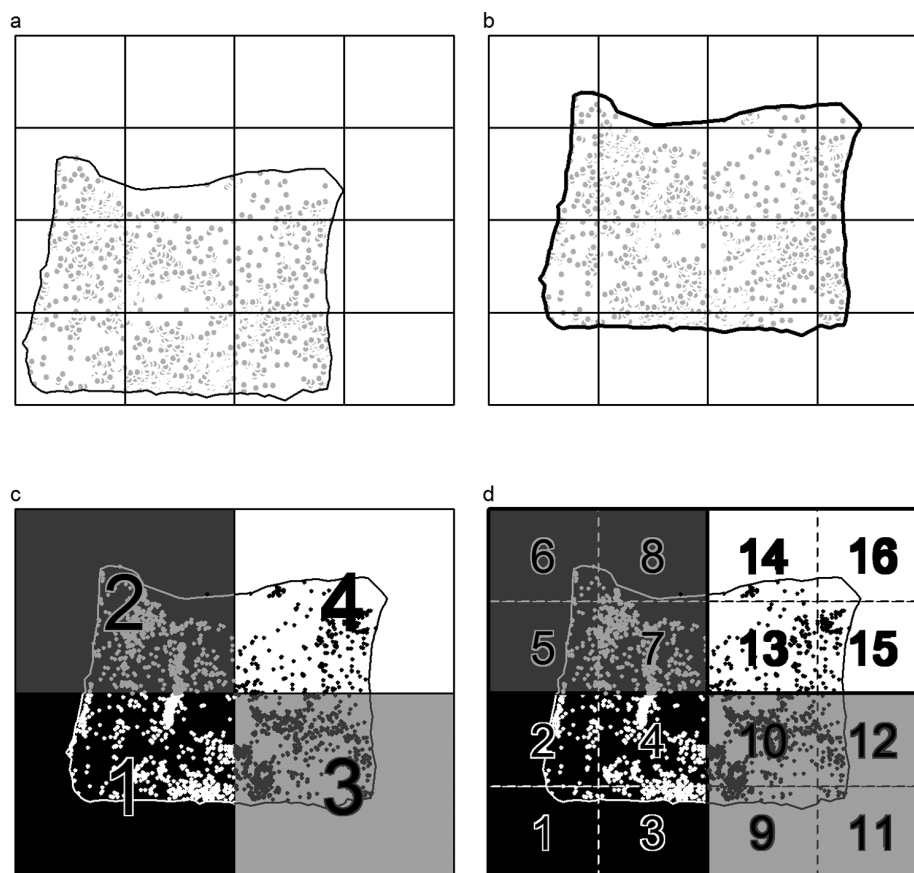


Figure 6.1 Sample frame of Oregon lake population. (a) Sample frame placed within systematic grid with extra row and column of empty grid cells with the grid size determined by the sample size and distribution of the elements to be sampled; (b) sample frame randomly shifted; (c) first-level hierarchical ordering of grid cells ($m = 1$); (d) second-level hierarchical ordering of grid cells ($m = 2$).

reason for requiring that any two lakes may occur in different cells is that this allows the possibility that both lakes could be in a sample that is selected. Technically, this ensures that their joint inclusion probability is greater than zero, a requirement for the Horvitz–Thompson variance estimator (Horvitz and Thompson 1952).

Next, the cells are numbered based on m hierarchical levels associated with the 2^m by 2^m grid of cells. At hierarchical level $m = 1$, the cells are numbered 1 to 4 and for illustration purposes colored from dark to light gray in Fig. 6.1c. At hierarchical level $m = 2$, each of the cells from the previous level ($m = 1$ in this case) are subdivided into four sub-cells. For simplicity the cells are numbered from 1 to 16 (Fig. 6.1d). Technically, a base 4 hierarchical numbering scheme is used to identify the cells and track their hierarchical location. Note that if another hierarchical level were required ($m = 3$), the number of cells would be $2^3 \times 2^3 = 64$. The hierarchical level is chosen

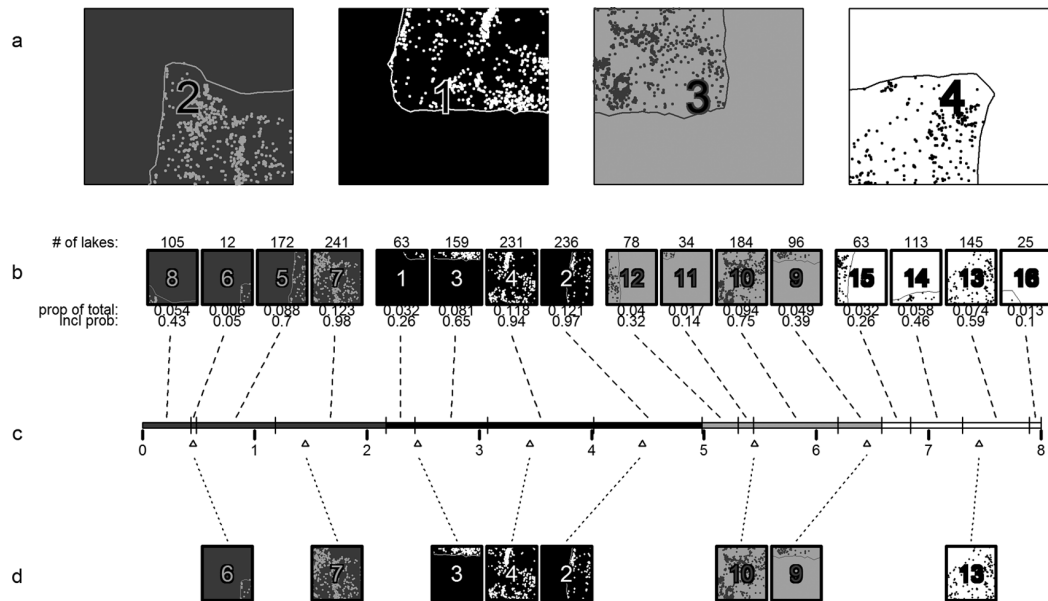


Figure 6.2 GRTS cell selection process. (a) Random ordering of first hierarchical level, $m = 1$. (b) Random ordering of second hierarchical level, $m = 2$, within first hierarchical level cells. (c, top) Creation of line with segments proportional to inclusion probabilities, p_i , and a total length of 8, which equals the desired sample size. (c, bottom) Selecting cells using a systematic sample with a random start of 0.45, and the triangles indicating the systematic set of points of 0.45, 1.45, 2.45, ..., 7.45. (d) The 8 cells selected based on their line segment containing one of the systematic set of points.

so that when the sample is selected it will include at most 1 lake from a cell. Thus the hierarchical level m depends on the sample size required and the spatial distribution of the lakes. The number of lakes that occur in a cell depends on the spatial distribution of lakes which impacts the probability of one or more lakes being selected in a cell. In our example with a sample size of $n = 8$ and the lake distribution given by the sample frame, the number of hierarchical levels required is $m = 2$, whereas a sample size of $n = 50$ using the same sample frame requires $m = 4$.

The next step is to complete a hierarchical randomization process. The process begins at hierarchical level $m = 1$ by randomly ordering the cell numbers and placing them in that order on a “line” (Fig. 6.2a). In this case the cell order is 2, 1, 3 and then 4. Then the next hierarchical level ($m = 2$) is randomly ordered within the previous hierarchical ordering. That is, the 4 sub-cells in each of the cells at hierarchical level $m = 1$ are randomly ordered independent of the other cells at that level (Fig. 6.2b). In this case, level $m = 1$ cell 2 has randomly ordered sub-cells 8, 6, 5, and 7; cell 1 has randomly ordered sub-cells 1, 3, 4, and 2; cell 3 has randomly ordered sub-cells 12, 11, 10, and 9; and cell 4 has randomly ordered sub-cells 15, 14, 13, and 16. These randomly ordered cells are placed on a “line” (Fig. 6.2c). Note that this hierarchical ordering process places the cells on the line so that portions of the state that were together at the first level remain

close together at the next level, as illustrated by the gray-scale color coding. If a third hierarchical level were required, the same process would be repeated for each of the 16 sub-cells in the second row and similarly for any additional hierarchical level. The basic idea of this hierarchical randomization process is to map two-dimensional space onto a line of cells where cells that are next to each other on the line tend to be next to each other in the original two-dimensional geographic space and at the same time introduce as much random ordering as possible. For an introduction to this concept of mapping a sample frame to a “line”, see Chapter 5, which illustrates this process for one-dimensional sampling.

For this example, assume for simplicity of illustration that we desire a sample size of 8 and that the spatial survey design specifies that the lakes should have an equal probability of being selected. This would be a simple random sample design if we were not interested in a spatially balanced GRTS design. Instead, we chose an equal probability GRTS design. How do we select 8 lakes from the 16 cells and make sure that they are spatially distributed across Oregon in the same way that the lakes in the sample frame are distributed across Oregon? First, we determine the number of lakes in each cell and their proportion of the total number of lakes in Oregon. Then we multiply these proportions by the desired sample size of 8 to obtain the inclusion probability of each cell. That is, the inclusion probability $\pi_i = n \frac{e_i}{E}$ for cell i where e_i is number of lakes (extent) in cell i , E is the total number of lakes in the sample frame (i.e. total extent) and n is the sample size. Note that the inclusion probability for the cell is the sum of the inclusion probabilities for all lakes, i.e. sample units, in the cell. In Fig. 6.2b all inclusion probabilities (“inc prob”) are < 1 . Each cell is mapped to a line segment of a length that is proportional to the inclusion probability of the cell. In this case, the cells are used to define segments of a line of length 8 (the desired sample size). Those cells with a small inclusion probability (fewer lakes) have shorter line segments and those with a large inclusion probability (more lakes) have longer line segments (Fig. 6.2c).

The next step is to select 8 cells from the line using a systematic sample process with a random start. Since the line has length equal to the sample size, we chose a random number between 0 and 1, say 0.45 to obtain a systematic set of points on the line at 0.45, 1.45, 2.45, . . . , 7.45 (indicated by triangles in Fig. 6.2c). Cells associated with the line segments with points from the systematic sample are selected (cells 6, 7, 3, 4, 2, 10, 9 and 13), Fig. 6.2d. Each of these cells will contribute a lake to the sample. The line segment associated with a selected cell is composed of smaller segments representing each lake. In this case because lakes are being selected with equal probability these segments are the same length (Fig. 6.3a). Individual lake line segments are randomly placed on the cell segment and the lake is selected where the systematic point occurs on the segment (at 0.45 or fifth lake in this example, Fig. 6.3b). This is completed for all cells to obtain the sample of 8 lakes (Fig. 6.3c). Table 6.1 summarizes the data and computations that are used to select the 8 cells and then to select the specific lake within each cell (assuming the lakes are randomly ordered with the cell). The combination of the hierarchical randomization process and the use of a systematic sample of the line based on it ensure that the sample will be spatially distributed across Oregon with

Table 6.1 GRTS selection of cells and lakes within cells for Oregon lake example.

Cell	# Lakes in cell	Proportion of lakes in cell	Cell inclusion probability	Cell line segment start	Cell line segment end	Systematic sample point	# Lake segments in cell to point	Lake selected within cell
8	105	0.054	0.429	0.000	0.429			
6	12	0.006	0.049	0.429	0.478	0.45	5.08	6
5	172	0.088	0.703	0.478	1.181			
7	241	0.123	0.985	1.181	2.167	1.45	65.71	66
1	63	0.032	0.258	2.167	2.424			
3	159	0.081	0.650	2.424	3.074	2.45	6.33	7
4	231	0.118	0.944	3.074	4.018	3.45	91.96	92
2	236	0.121	0.965	4.018	4.983	4.45	105.58	106
12	78	0.040	0.319	4.983	5.302			
11	34	0.017	0.139	5.302	5.441			
10	184	0.094	0.752	5.441	6.193	5.45	2.21	3
9	96	0.049	0.392	6.193	6.586	6.45	62.83	63
15	63	0.032	0.258	6.586	6.843			
14	113	0.058	0.462	6.843	7.305			
13	145	0.074	0.593	7.305	7.898	7.45	35.46	36
16	25	0.013	0.102	7.898	8.000			
Total	1957	1.000	8.000					

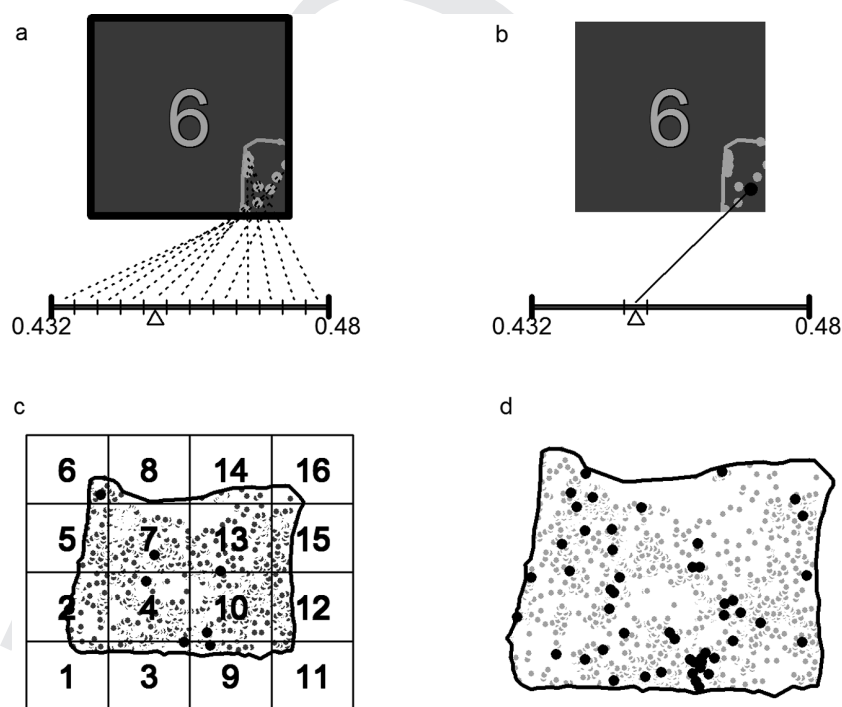


Figure 6.3 Illustration of the random selection of a lake within a selected cell. (a) random ordering of lakes on line with equal length assigned to each lake. (b) Selection of lake when random location of point occurs within the lake's line segment. (c) GRTS sample of size eight. (d) GRTS sample of size 50 where number of hierarchical levels determined so that total inclusion probability of each grid cell was <1 .

a spatial distribution that is similar to the spatial distribution of lakes. This is better illustrated in Fig. 6.3d with a sample size of 50. Note that even though two lakes are close together, it is possible for both of them be in a sample due to the random shift of the sample frame (see lakes selected in cells 3 and 9). Also the selection process within a cell is the same as selecting a simple random sample of size 1 from the lakes in the cell.

How is the hierarchical level, i.e. “ m ”, determined? The value for m is selected so that all $2^m \times 2^m$ grid cells in the hierarchy have an inclusion probability π less than 1. The inclusion probability depends on the total sample size, the number of lakes within the cell and the total number of lakes in the population. Because the grid is randomly located, the inclusion probability can only be calculated after the randomization, as the number of lakes in a cell is unknown until then. If we assume the lakes are uniformly distributed in space and the desired sample size is n , then an initial value for m is the next larger integer of $\log_4(n)$. If n equals 100, then m would be equal to 4. When the lakes are clustered, then m will mostly likely need to be greater to ensure the inclusion probability is less than 1. The determination of m is automatically calculated in the “*spsurvey*” *grts()* function so the user does not need to specify m .

Basic GRTS algorithm for linear networks and polygons

We used lakes as a collection of points in space to illustrate the GRTS algorithm. The same basic algorithm is used for a linear network, e.g. streams, roads or trails, and for polygons, e.g. forests, estuaries or lakes sampled as polygons. The process is the same through the steps illustrated in Fig. 6.1 and Fig. 6.2a,b except that instead of points in cells we now have linear segments or polygons in cells. In the step illustrated by Fig. 6.2c, the algorithm differs in the determination of the line segment length associated with each cell and then is the same for the step illustrated by Fig. 6.2d. The subsequent step of selecting a sample point on the linear network within a cell or selecting a sample point within a polygon in a cell differs from that illustrated for lakes. Note that linear networks, such as stream networks, are digitally represented by straight line segments of different lengths. While the sample unit is a point for discrete natural resources, i.e. a lake in the example, the sample unit for linear networks or polygons is any location on the linear network or any point within the polygons. In both cases the number of sample units is infinite as opposed to finite for a point natural resource.

What differs in the step illustrated by Fig. 6.2c, i.e. in constructing the line? In the lake example, we determined the number of lakes in each cell. We call this the extent of the sample frame within the cell. For a linear network the extent is the total length of the segments that make up the linear network within the cell. If a segment extends outside the cell it is clipped at the cell boundary and only that portion of the segment length is used. The rest of the segment is included in the extent computations of the neighboring cell. The inclusion probability is $\pi_i = n \frac{e_i}{E}$ for cell i , where e_i is extent (linear network length) in cell i , E is the total extent (total linear network length in the sample frame), and n is the sample size. For an areal sample frame (polygons), the extent is the total area

of the polygons that are within the cell. If a polygon extends outside the cell it is clipped at the cell boundary and only that portion of the polygon area within the cell is used. Extent is in units of numbers for a point sample frame, length for linear network sample frame and area for a polygon sample frame. The line segments associated with each cell for a linear network or areal sample frame have length equal to the extent within each cell. Otherwise the construction of the line is the same as for points as is the selection of the systematic sample of points on the line.

How is a random point selected on the linear network within a cell? These segments within a cell are randomly ordered and then placed on the line in that order with length equal to their inclusion probability (which is proportional to their length). Using the example based on lakes, the segment on the linear network selected from the first cell would be the line segment that included 0.45. For illustration assume it is the first line-segment and that it has an inclusion probability of 0.03. The cell (6) begins its line segments at 0.429 and ends at 0.478 and the first linear-network segment begins as 0.429 and ends at 0.46. The selected point would be $100 \cdot (0.45 - 0.429) / (0.46 - 0.429) = 66.7\%$ up the line segment from the beginning of the line segment. Given we know the percent we can determine its location on the linear-network segment and then its x -, y -coordinates. Note that sample points may be selected anywhere on the linear network.

How is a random point selected in the polygons within a cell? Areal sample frames are collections of polygons of different areas. The polygons within a cell are randomly ordered and then placed on the line in that order with length equal to their inclusion probability (which is proportional to their area). Using the example based on lakes, the polygon selected would be the polygon that included 0.45, assume it is the first polygon and that it has an inclusion probability of 0.03. The cell (6) begins its line at 0.429 and ends at 0.478 and the line segment associated with the first polygon begins at 0.43 and ends at 0.46, consequently the first polygon would be selected. Then within that polygon a random point would be selected. The latter is accomplished by selecting random x and y coordinates while ensuring the coordinates occurred within the polygon. While it may seem simpler to just select random x - and y -coordinates within the entire selected cell and ensuring the random point is in a polygon, this has drawbacks. First, the cell may include only one polygon that is very small compared to the size of the cell. Picking a random x , y coordinate pair based on the cell will have a low probability of having the point fall within the small polygon requiring many random draws before the point falls within the cell. Second, more complex designs may require some polygons to have a probability of being selected that is not proportional to their area (unequal probability designs discussed in next section).

GRTS applied to stratified and unequal probability spatial survey designs

Environmental monitoring programs typically have monitoring objectives or other requirements that require the use of complex survey designs. Examples of complex designs applied to aquatic resources are the US EPA National Aquatic Resource

Surveys (see <http://water.epa.gov/type/watersheds/monitoring/nationalsurveys.cfm>) for lakes (points), streams and rivers (linear network), coastal waters (polygons), and wetlands (polygons). Four basic design options (reviewed in Chapter 5) are simple random sample, stratified random sample, unequal probability sample and stratified unequal probability sample designs. Each of these is also an option with GRTS designs. We have already described the GRTS equal probability spatial survey design. It is equivalent to a simple random sample except that it ensures each sample realization has a spatial distribution that is similar to the spatial distribution of the environmental resource. Next we discuss stratification and then unequal probability sampling. The fourth basic design, a stratified unequal probability spatial survey design, simply combines stratification with an unequal probability design.

A stratified GRTS spatial survey design is a simple extension of the basic GRTS algorithm. Stratification divides the sample frame into independent sample frames that collectively equal the entire sample frame. Since each stratum has its sample selected independently of the other strata, the GRTS algorithm is simply applied to each stratum. For example, see Appendix 6.1 for *grts()* commands used if the population of Oregon lakes was stratified into “mountain” and “xeric” region groups.

Unequal probability survey designs are another basic survey design and can be used as an alternative to stratified survey designs when the objective is to reduce variance of estimates (Lohr 1999). An unequal probability GRTS spatial survey design involves a simple modification to the basic GRTS algorithm to incorporate the unequal probability information. The basic GRTS algorithm requires that inclusion probabilities be determined for each cell. For the lake example these are proportional to the number of lakes in a cell since each lake was to be selected with equal probability. Suppose that instead of placing lakes in mountain vs. xeric areas in two different strata, we wanted to select “mountain” and “xeric” lakes with unequal probability, say a xeric lake with twice the probability of a mountain lake. Then when calculating the inclusion probability we would assign a line-segment length value of 1 to all mountain lakes and a line-segment length value of 2 to all xeric lakes. Then we would sum those values for all lakes in a cell to obtain the cell extent e_i and the total extent E . Note that now E is no longer equal to the number of lakes, but is equal to the number of mountain lakes plus twice the number of xeric lakes. Otherwise the selection process remains the same. For linear networks, the process is similar except that a segment length is multiplied by the values to get the inclusion probabilities. For polygons, the polygon areas are multiplied by the values to get the inclusion probabilities. Typically a user will know the number of samples desired for mountain and xeric lakes. For example, the user may want the expected sample size to be 50 for mountain lakes and 50 for xeric lakes for a total sample size of 100. (The “*spsurvey*” *grts()* algorithm uses the sample size information for the two categories of lakes to determine the appropriate values to use when calculating the inclusion probabilities; see Appendix 6.1.) Note that in contrast to a stratified sample the unequal probability sample does not guarantee exactly 50 lakes in each category – only that on average over repeated sampled draws 50 lakes will be in each category with the total always being 100 lakes. This is an inherent feature of unequal probability sampling.

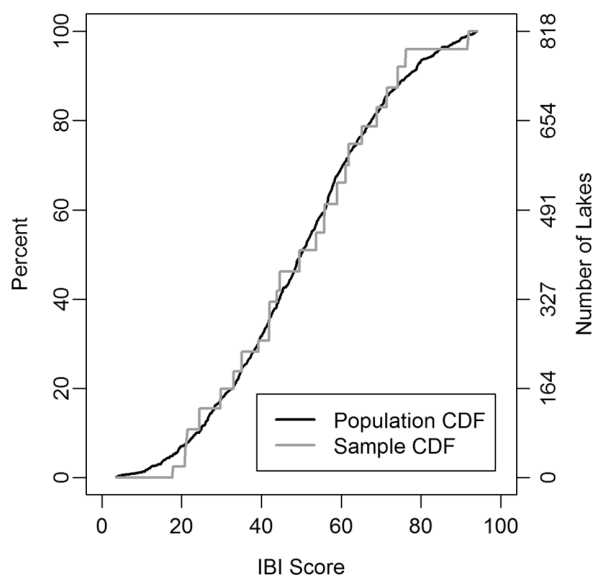


Figure 6.4 Illustrative cumulative distributions (CDFs) for a simulated index of biotic integrity (IBI) score variable for a target population of lakes and a miniature of the population (sample) using a sample of size 25. See Chapter 14 for additional discussion of CDFs.

Spatial balance and spatial balance metric

Earlier we stated that monitoring designs should be based on a representative sampling process that results in a representative sample that is a miniature of the target population and that one way to achieve this was to use a process that resulted in the sample having the same spatial distribution as the population. We then described GRTS spatial survey designs and claimed that they satisfy these requirements by selecting spatially balanced samples. What do we mean by “a sample having the same spatial distribution as the population”? In this section, we first define “spatial balance” to give specific meaning to this phrase and then we define specific metrics to determine if a specific realization (a sample) is spatially balanced.

First, we give a quantitative definition measuring what it means for a representative sample to be a miniature of the population. Assume a single variable, x , is of interest. The population is completely characterized by its cumulative distribution function $F(x)$. A sample from the population can be used to construct a sample distribution function $\hat{F}_n(x)$ (Fig. 6.4; Chapter 14). One possible metric to measure representativeness can be based on the Kolmogorov–Smirnov statistic $D_n = \max |\hat{F}_n(x) - F(x)|$. This metric measures the maximum deviation of the sample distribution function from the population distribution function. When D_n is small the sample empirical distribution function is “close” to the population cumulative distribution function. Consequently, inferences based on the sample empirical distribution function will be the same as for the population. Most monitoring programs measure many indicators so that this condition would need to



Figure 6.5 Dirichlet tessellation for a sample of size 25 from Oregon lake population.

be satisfied for all indicators and preferably their multivariate empirical distribution function. While this metric defines what we mean by a representative sample being a miniature of the population, it is not useful in practice since the population distribution function is unknown.

Given the difficulty in determining whether a sample is a miniature of the population, instead we focus on metrics that measure the spatial balance of an equal probability, non-stratified sample of size n . Our spatial balance metrics are based on the Dirichlet tessellation (also called Voronoi or Thiessen polygons) associated with the sample of size n . The Dirichlet tessellation is the set of polygons over the study region created such that all locations within any given polygon are closer to one of the sample points than to any other sample point (Fig. 6.5). These polygons are then used to calculate the relative extent (i.e. the number, length or area) or proportion p_i of the target population that is within each polygon associated with site i . If the sample is spatially balanced we would expect that the proportion of the target population in each polygon would be $1/n$. Figure 6.5 illustrates the Dirichlet tessellation for a sample of size 25 for the Oregon lake population.

Based on the proportions p_i from the Dirichlet tessellation of the sample, we define three alternative spatial balance metrics of a selected sample:

- (i) The J_p spatial balance metric is Pielou's evenness index, which is based on the Shannon–Wiener index. It is defined as $J_p = -\sum_i^n p_i \ln p_i / \ln n = H/H_{\max}$, where H is the Shannon–Wiener diversity index and H_{\max} is its maximum value (Legendre and Legendre 1998). This is one of several diversity indices used to measure diversity in categorical data. The maximum for H , H_{\max} , occurs when all p_i are equal resulting in the maximum value being equal to $\ln n$. We define spatial balance using Pielou's evenness index where perfect spatial balance has $J_p = 1$, i.e. perfect evenness of the target population associated with all sample sites.

- (ii) The X_p^2 spatial balance metric is equal to $X_p^2 = \sum_{i=1}^n (p_i - 1/n)^2 / (1/n)$. This metric is simply the common chi-square statistic $\sum (O - E)^2 / E$. Perfect spatial balance occurs when $X_p^2 = 0$.
- (iii) The S_p spatial balance metric is defined as $S_p = \sqrt{\sum_{i=1}^n (p_i - \bar{p})^2 / (n - 1)}$, the standard deviation of proportions. In this case, perfect spatial balance occurs when $S_p = 0$.

We conducted a simulation study on the spatial-balance properties for a GRTS equal probability survey design vs. those of a simple random sample (IRS) survey design that does not use location as part of the sample selection process. Five alternative target populations were used: (a) lakes in Oregon as a point population using lake centroid as the point; (b) hydrologic units for Oregon (see <http://www.ncgc.nrcs.usda.gov/products/datasets/watershed/> for definition of hydrologic units) as a point population using unit centroid as the point; (c) perennial streams in Oregon as a linear network population; (d) hydrologic units in Oregon as polygon population; and (e) lakes in Oregon as polygon population. For each target population 1000 samples of size 5, 10, 25, 50 and 100 were selected using GRTS and IRS designs, i.e. $5 \times 5 \times 2 \times 1000 = 50000$ samples. For each sample we computed the Dirichlet tessellation and calculated the value of the three spatial balance metrics. In this chapter, we only report the results for Oregon lake population treated as points to compare the spatial balance properties of GRTS and IRS designs. The results for the other populations were similar.

The spatial balance properties for the GRTS and IRS simulated samples are presented in Fig. 6.6. The GRTS evenness metric median is greater than the IRS median for all sample sizes as are the 75th percentile and the maximum (Fig. 6.6a, left). For sample sizes of 25, 50 and 100 the evenness metric values for GRTS samples are almost entirely greater than 75% of the values for IRS samples. In all cases, larger values of the evenness metric reflect better spatial balance. For sample sizes of 5 and 10, the GRTS evenness distribution is better in terms of spatial balance than for IRS. For both GRTS and IRS samples the variability of the evenness metric increases as the sample sizes become smaller, although the increase is greater for IRS than for GRTS. Overall, the simulation shows that any single realization from a GRTS spatial survey design is more likely to result in a sample that is spatially balanced than any single realization from an IRS spatial survey design. The long tail for IRS compared to that for GRTS emphasizes the relative poor performance of IRS in achieving spatial balance compared to GRTS. The chi-square and standard deviation metric simulation results are consistent with those from the evenness metric, although small values now indicate better spatial balance (Fig. 6.6a, center and right). The relationships among the 3 metrics for a sample size of 25 are shown for both GRTS (Fig. 6.6b) and IRS (Fig. 6.6c). The t_3 metrics are highly correlated and have a nearly linear relationship, but we present all 3 metrics for consideration by readers who may want to implement a similar examination of comparative spatial balance. Variability increases among the metrics as spatial balance decreases, although this is less so for the chi-square and standard deviation metrics.

Figure 6.7 shows the most and least spatially balanced sample realizations from the 1000 simulations for GRTS and IRS, based on the evenness metric. It is visually difficult

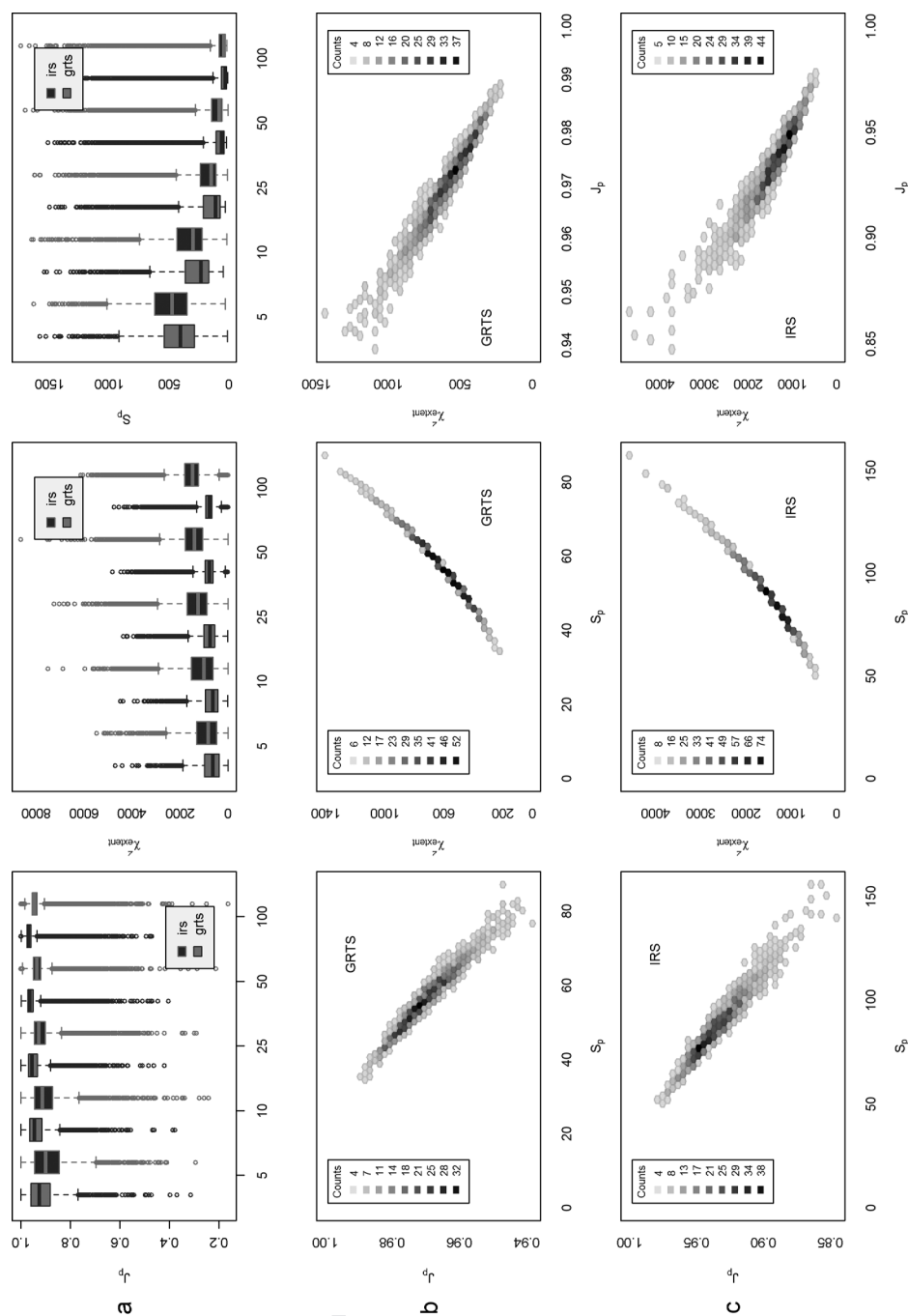


Figure 6.6 Spatial balance metrics for Oregon lake population based on 1000 simulations for GRTS and IRS (simple random) samples of size $n = 5, 10, 25, 50$, and 100 . (a) Results for GRTS vs. IRS for (from left to right) evenness metric, chi-square metric, and standard deviation metric. (b) Pairwise relationships among metrics for GRTS, $n = 25$. (c) Pairwise relationships among metrics for IRS, $n = 25$.

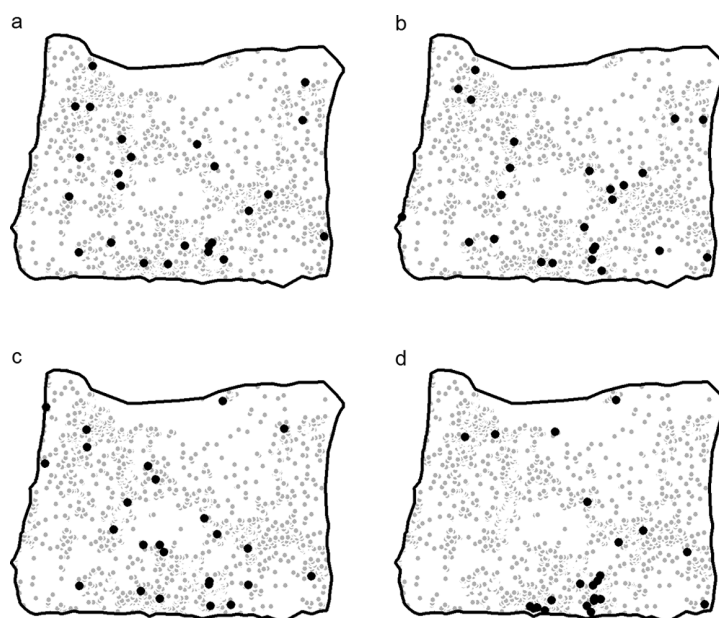


Figure 6.7 For sample size of 25, realization of a sample from 1000 simulated samples that is (a) the most spatially balanced GRTS sample ($J_p = 0.99$); (b) the most spatially balanced IRS (simple random) sample ($J_p = 0.98$); (c) the least spatially balanced GRTS sample ($J_p = 0.93$); (d) the least spatially balanced IRS sample ($J_p = 0.85$).

to distinguish the most (a) and least (c) spatially balanced samples for GRTS. We would expect this to be the case since the evenness metric values are very similar, $J_p = 0.99$ and $J_p = 0.93$, respectively. The most spatially balanced IRS sample is similar visually to the GRTS most and least spatially balanced samples; its evenness metric value is 0.98. The IRS least spatially balanced sample ($J_p = 0.85$) is noticeably different with the sample not reflecting the spatial distribution of the Oregon lakes.

Overall, Figs. 6.6 and 6.7 illustrate the difference in spatial balance between GRTS and IRS. Both can produce a sample that has good spatial balance. The spatially balanced samples that are possible using the GRTS representative sampling process are a subset of all the samples that are possible using the IRS representative sampling process. Consequently, the IRS samples with low spatial balance, as indicated by the J_p evenness metric, do not occur for GRTS.

Additional issues

Spatial balance properties when entire sample is not used

A spatial survey design implemented as part of an environmental monitoring program typically encounters anticipated as well as unanticipated problems. For example, the sample locations from the design may be (i) determined to not be a member of the

target population; (ii) found to be inaccessible as a result of a landowner denying access; (iii) found to be physically inaccessible due to safety concerns or excessive cost to access; or (iv) not sampled due to insufficient time to complete the study within the allowed time or as a result of a reduction in funds. Each of these problems leads to a smaller sample size than what was originally planned. The smaller sample size resulting from the first three problems can be addressed by constructing the design so that additional sample locations are available if needed. As an example, if the study planned to collect data from 50 sample locations and only 40 of the 50 selected sample locations could be sampled, then if additional sample locations were available they could be used to obtain data from 10 replacement sample locations. The objective is simply to have the study result in data for the planned sample size of n , i.e. 50 sample locations in the example. Meeting the desired sample size does not address other issues arising from the “non-response” of locations not sampled (such as bias; see discussion in Chapters 3, 5). The question is whether it is possible to augment the original “base” survey design with a supplemental design of “over sample” sample locations.

One approach is to implement a spatial survey design where the sample size is the total of the desired sample size n (base sample) plus an additional sample size o (over sample) that would only be used if necessary. This can be done for an IRS spatial survey design, simply by selecting a sample of size $n + o$, using the first n sample locations first, and then using the o locations in the order they appear in the randomized list from the design. For a GRTS spatial survey design, the sample locations selected will be ordered according to the spatial hierarchical order use in the GRTS algorithm. Heuristically, if $n = 50$ and $o = 25$, then the first 50 sample locations will be located in only 75% of the study region. Stevens and Olsen (2004) define a reverse hierarchical ordering process that addresses this problem. It is not part of the basic GRTS spatial survey design process and is only relevant when the entire sample size is not used (e.g. Olsen and Peck 2008, Olsen *et al.* 2009). Note that the “spsurvey” *grts()* function demonstrated in Appendix 6.1 automatically does the reverse hierarchical ordering of the selected sample.

We investigated the spatial balance properties of IRS and GRTS spatial survey designs when only a portion of the entire sample ($n + o$ from above) was used. This addresses the question of whether the reverse hierarchical ordering algorithm results in a spatially balanced sample when the entire sample is not used. For each simulated sample, the sample points were used in the order they appeared in the sample. For GRTS simulated samples this was after the reverse hierarchical ordering was applied. Sample points were added 1 point at a time up to the maximum sample sizes of 50 or 100. The results for Oregon lakes are given in Fig. 6.8. For any specific subsample size, the summary is similar to the boxplot summaries in Fig. 6.6a, except that here the minimum, 25th percentile, median, 50th percentile and maximum values are plotted. When the subsample size is close to the planned sample size, the spatial balance properties are similar to the full sample size of 50 (Fig. 6.8a) and of 100 (Fig. 6.8b). As the subsample size becomes smaller, the spatial balance properties of both IRS and GRTS spatial survey designs deteriorates, although GRTS continues to be better than IRS until very small sample sizes where it is similar to IRS.

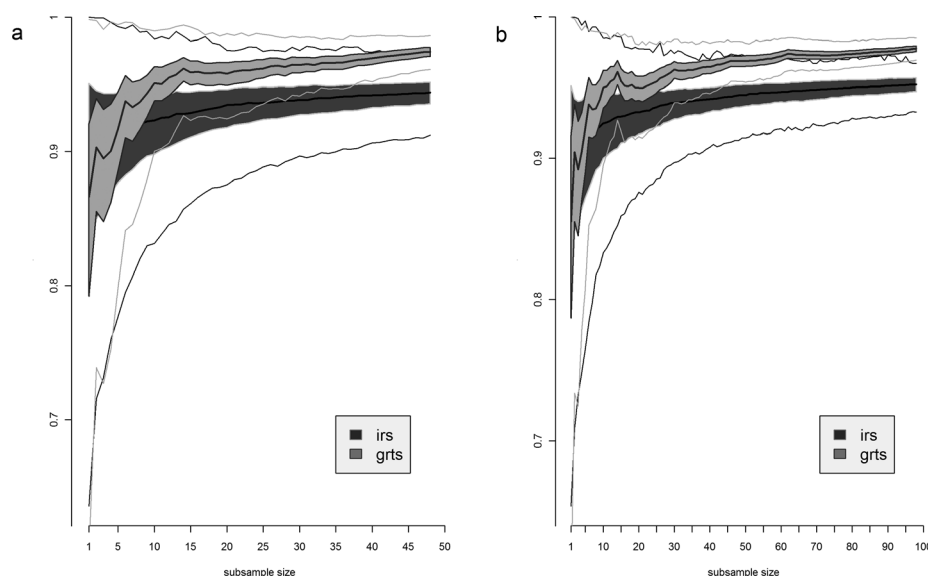


Figure 6.8 Spatial balance properties as a function of sample size for GRTS and IRS (simple random) samples based on 1000 simulations of the Oregon lake population. The y -axis is the evenness spatial balance metric (J_P). Sample points were added one point at a time up to the maximum sample sizes of (a) 50 or (b) 100. GRTS samples were arranged in reverse hierarchical order (see text). The shaded area includes the middle 50% of the simulated metric values (25th to 75th percentiles), dark line is the median (50th percentile) and the outer lines the minimum and maximum metric values.

GRTS applied to temporal designs

Chapter 7 discusses temporal (revisit) designs that are useful for trend detection. The most direct use of a spatial survey design for detecting trends is simply to select a spatially balanced sample, collect measurements on indicators of interest in the first time period, and then revisit the same sites in subsequent time periods. This is commonly called an annual revisit (or always revisit) panel design. (A panel is a collection of sites that have the same revisit pattern over time.) Another trend design is to never revisit the same sites in subsequent time periods. For example, if the monitoring program is planned for 20 years, then each year a new panel of sites is visited for a total of 20 panels. Other panel designs are discussed in Chapter 7.

Regardless of the panel design, a spatially balanced sample for trend detection can be selected as follows. Assume each panel will have n_i sites for $i = 1$ to P panels. The total number of unique sites is $n = \sum_{i=1}^P n_i$. A spatially balanced sample of n sites is selected, and the reverse hierarchical ordering algorithm is used. Then, the first n_1 are assigned to panel 1, the next n_2 are assigned to panel 2, and so on. Each panel of sites is a spatially balanced sample. Note that each panel is a spatially balanced sample for the population so that status estimates are possible for each time period as well as trend and change estimates across and between time periods.

Future research and development

Our description of spatial balance focused on spatial balance with respect to the target population. If the population is clustered, then the spatially balanced sample will be clustered. In some situations population sample units that are close together in space may be expected to be similar in their characteristics, i.e. they are correlated. It may be desirable to have a different type of spatial balance in this case – spatial balance with respect to geography. In the Oregon lake example, this means that the sample would be more uniformly spread across the state. Geographic spatial balance can be defined, although it is more difficult to implement. The “*spsurvey*” *grts()* algorithm can be applied but the user must create the inclusion probabilities based on a two-dimensional spatial density estimation of the population. Further research is required before this will be available.

Summary

In this chapter we have focused on spatial survey designs, and specifically spatially balanced survey designs and GRTS, and their use in monitoring programs for natural resources. Natural resources occur in two-dimensional space and can be modeled as points, linear networks, and polygons. Building on the concept of a representative sample as a miniature of the population, we described how spatially balanced survey designs can achieve that representativeness. We also distinguished a representative sampling process from a representative sample, noting that a representative sampling process does not necessarily result in a representative sample for every sample realization. We then defined spatial balance metrics and showed how a GRTS spatial survey design results in better spatial balance than a simple random sample (IRS). As GRTS spatial survey designs can also include stratification and unequal probability of selection, our conclusion is that GRTS spatial survey designs should be used for all natural resource monitoring programs where a probability survey design is the appropriate design to meet the monitoring objectives.

Acknowledgments

Marc Weber provided us with the GIS shapefiles used in our simulations. We appreciate reviews provided by Phil Larsen and two anonymous referees. The information in this document has been funded by the US Environmental Protection Agency (USEPA). This manuscript has been subjected to review by the National Health and Environmental Effects Research Laboratory’s Western Ecology Division and approved for publication. Approval does not signify that the contents reflect the views of the Agency, nor does mention of trade names or commercial products constitute endorsement or recommendation for use.

Appendix 6.1. Implementing GRTS spatial survey designs in R

The R package “spsurvey” implements the GRTS spatial survey design algorithm and includes vignettes illustrating how the function *grts()* within the package is used to select a spatially balanced sample given a design specified by the user. The “spsurvey” package is available from the R software website (<http://www.r-project.org/>) or the aquatic resource monitoring web site (<http://www.epa.gov/nheerl/arm/>). The function *grts()* selects samples from sample frames based on points, linear networks or polygons to obtain a GRTS spatial survey design. In addition to equal probability GRTS spatial survey designs, the function *grts()* includes options for stratified GRTS spatial survey designs, unequal probability GRTS spatial survey designs based on user specified categories, or unequal probability GRTS spatial survey designs where the unequal probabilities are proportional to a user-specified continuous variable.

The following is an example of using the R package “spsurvey” to select the sample of 50 lakes. Three basic steps are required: (i) reading in the sample frame, (ii) specifying the design, and (iii) selecting the sample:

```
att <- read.dbf('Oregon_Lakes')
dsgn <- list(None=list(panel=c(Base=50), seltype= 'Equal'))
lakes <- grts(design=dsgn,
              DesignID='ORLakes',
              type.frame='finite',
              src.frame='shapefile',
              in.shape='Oregon_Lakes',
              att.frame=att,
              shapefile=TRUE,
              out.shape="OR_Lake_Design_Sites",
              prj="Oregon_Lakes")
```

The first line reads in the attribute file that is part of the ESRI shapefile of Oregon lakes. It contains attributes associated with each lake, such as name, area, etc. The second line specifies the GRTS spatial survey design requirements. In this case “None” means the design is not stratified, “Base = 50” means that the sample size is 50 and all are assigned to a panel named “Base”, and “seltype = ‘Equal’” means that the lakes should be selected with equal probability. The *grts()* function provides the information the algorithm requires to select the lakes using the design “dsgn”. “DesignID” specifies a prefix that will be added to the sample number assigned to the lakes and used as a “siteID”. The term “type.frame” indicates that the sample frame is a collection of points (as opposed to a linear network or set of polygons); “src.frame” specifies that the source of the sample frame is an ESRI shapefile and “in.shape” gives the name of the ESRI shapefile. The term “att.frame” specifies that the attribute information is contained in “att”, an R data.frame; “shapefile = TRUE” indicates that the selected sample should be written out as an ESRI point shapefile and “out.shape” gives the name to use for that shapefile. The term “prj = ‘Oregon_Lakes’” tells *grts* that the map projection of the

output shapefile is the same as the input ESRI shapefile. The R object “lakes” contains the same information as the output shapefile and can be used within R.

If the lakes were stratified into mountain and xeric strata with sample size of 50 in each, this would be specified in a design requirement statement as follows:

```
dsgn <- list(mountain=list(panel=c(Base=50), seltype= ‘Equal’),  
xeric=list(panel=c(Base=50), seltype= ‘Equal’))
```

The only change to the *grts()* function would be to add a line with “stratum = ‘stratum’” where ‘stratum’ would need to be the name of one of the columns in the attribute portion of the shapefile and each lake would have a stratum attribute of either “mountain” or “xeric”. Note that exactly 50 lakes would be selected in each of the two strata.

For an unequal probability sample of lakes where the expected sample size is 50 for mountain lakes and 50 for xeric lakes, the design statement used in *spsurvey* would be:

```
dsgn <- list(None=list(panel=c(Base=100), seltype= ‘Unequal’),  
caty.n=c(mountain=50, xeric=50))
```

where “seltype” now specifies that unequal probability of selection is to be used to select the sample of size 100 and that the inclusion probabilities should be constructed within *grts()* so that the expected sample size is 50 for mountain lakes and 50 for xeric lakes. The *grts()* algorithm uses the sample sizes and the sample frame information for the two categories of lakes to determine the appropriate values to use when calculating the inclusion probabilities. The *grts()* statement requires the addition of “mdcaty = ‘laketype’” where ‘laketype’ is an attribute in the shapefile with each lake having an attribute of either “mountain” or “xeric”.