

MBHdesign: an R-package for efficient spatial survey designs

Scott D. Foster*

26 Nov 2020

Abstract

1. A considered survey design will generate data that is representative of the population that the sample is taken from. All good design takes is a little thought, some information and some good software tools.
2. Spatially-balanced randomisation with unequal inclusion probabilities is a modern and efficient design method. These designs are embedded within sampling theory and should be easy to generate.
3. The R-package **MBHdesign** allows field researchers easy access to these designs. It implements point-based and transect-based methods and allows for the incorporation of legacy sites.
4. The functionality of the package is illustrated with example designs in an environmentally heterogeneous region.

1 Introduction

Robust science can only be achieved using a rigorous and carefully planned scientific process (e.g. Leek & Peng 2015; Hayes *et al.* 2019). A pivotal link in this process is survey design, which transforms the research questions into a formal plan about how data is to be collected. It is important that the survey plan is designed so that the resulting data are: (i) representative of the population under investigation so that inference is valid; and (ii) information rich so that uncertainty, about inferences for the research question, is reduced. Without these attributes, surveys are less likely to deliver fit-for-purpose data (Hayes *et al.* 2019). To aid good survey design, the **MBHdesign** R-package is presented.

*Data61 CSIRO, Hobart, TAS, Australia, scott.foster@data61.csiro.au

21 The `MBHdesign` package leverages off the modern design strategy of spatial-balance (see Stevens & Olsen
22 2004) using Balanced Acceptance Sampling (BAS Robertson *et al.* 2013) with unequal inclusion probabilities
23 (see Thompson 2012). Spatially-balanced designs are robust because they are based on randomisation, which
24 guards against unintentional bias (see Altman 1991). The BAS method uses quasi-random numbers, which
25 in the context of survey design can be treated as random. Spatial-balance increases efficiency in two ways,
26 compared to simple randomisation. Firstly, it approximately balances over all spatially-smooth covariates
27 that are not considered or even measured (Grafström & Lundström 2013). Secondly, when considering
28 model-based analysis methods spatial-balance reduces the spatial autocorrelation between observations.
29 Unequal inclusion probabilities increase efficiency by allowing researchers to increase sampling effort for
30 certain environmental conditions, such as those that are likely to have higher variance (e.g. Thompson 2012).
31 In ecology, where variance often increases with the mean, sites with higher abundance should be sampled
32 more often. One consequence of unequal inclusion probabilities is that unweighted means will no longer be
33 an unbiased description of the sample.

34 In addition to providing efficient BAS designs, `MBHdesign` makes accessible two extensions. Firstly, `MBHdesign`
35 can accommodate *legacy sites*, which are often sites that have a long historical time-series and continuation
36 of this time-series is beneficial. The legacy sites are incorporated using the methods described in Foster *et al.*
37 (2017) and incorporates the legacy sites' locations into the spatially-balanced sample. Secondly, `MBHdesign`
38 can sample units that are *transects* rather than a *points*. The method described in Foster *et al.* (2020) is
39 implemented, which attempts to spatially-balance the centres of the transects whilst simultaneously respecting
40 the inclusion probabilities for each point. Transect sampling is only available in `MBHdesign`.

41 Several packages exist for generating spatially-balanced designs (Kermorvant *et al.* 2019). These include:
42 `spsurvey` (Kincaid *et al.* 2019) that implements the generalized random-tessellation stratified (GRTS)
43 algorithm (Stevens & Olsen 2004); `SDraw` (McDonald & McDonald 2020) that implements a range of spatially-
44 balanced methods including BAS, and; `BalancedSampling` (Grafström & Lisic 2019) that implements the local
45 pivotal method (LPM Grafström 2012) and spatially correlated Poisson sampling (SCPS Grafström *et al.* 2012).
46 All these packages provide a good platform for generating spatially-balanced designs, but they have different
47 foci in terms of algorithms, functionality, scope, computation requirements and user-interface. `MBHdesign` is
48 distinguished from these packages in that it provides field scientists with a tool that has a simplified yet flexible
49 interface to generating designs that are: 1) superior spatial-balanced (unlike `spsurvey`), 2) based on unequal
50 inclusion probabilities (unlike `SDraw`), and 3) based on computationally efficient methods that scale well to
51 large problems (unlike `BalancedSampling`). In addition `MBHdesign` provides functionality for incorporating
52 legacy sites (Foster *et al.* 2017) into the new survey design, as well as generating designs for transect-based

53 sampling platforms (Foster *et al.* 2020, but with extra computational demand). The package is freely available
 54 from the Comprehensive R Archive Network (CRAN), at <https://cran.r-project.org/package=MBHdesign>
 55 with a GNU GPL-3 license. This document was created with `MBHdesign` version 2.1.8.

56 2 The `MBHdesign` Package

57 The `MBHdesign` R-package is purposefully simple, with only three main functions and three ‘helper’ functions
 58 (see Table 1 for an overview). All functions have only a small number of arguments that are mandatory, but
 59 finer control can be achieved by changing default values.

Table 1: Functions available in `MBHdesign`. The first three functions are the primary functions in the package.

Function	Description
<code>quasiSamp</code>	Generate a BAS sample in arbitrary dimensions.
<code>alterInclProbs</code>	Adjusts inclusion probabilities to respect the locations of legacy sites.
<code>transectSamp</code>	Generates a spatially-balanced design for transect-based sampling.
<code>findDescendingTrans</code>	Finds transects within the survey area that run down gradients.
<code>findTransFromPoint</code>	Finds transects within survey area that originate from a given set of points.
<code>modEsti</code>	Simple model estimation method introduced in Foster <i>et al.</i> (2017).

60 `MBHdesign` takes a consistent and simple approach to spatial data: all functions can be called using a dense
 61 grid of points stored in a data frame (a non-spatial object). The data frame can be constructed using the
 62 `raster` package (Hijmans 2018), see supporting R-scripts for example code. It is recommended that the grid
 63 contains all locations within the extent of the area to be surveyed, and that an equal-area projection is used
 64 to maximise the efficacy of spatial balance. That is a (hyper-)rectangle that contains all the possible sampling
 65 locations. Those locations that are not to be sampled should be encoded as `NA`. For point-based designs,
 66 these locations can be deleted from the data set, but they are required for transect-based designs.

67 2.1 Hippolyte Rocks Examples

68 The functions within `MBHdesign` will be illustrated by way of creating designs for the marine environment
 69 surrounding the islands of the Hippolyte Rocks, Australia. For illustrative purposes, the bathymetric data

70 used here are a spatially-degraded version of those described in Nichol *et al.* (2009) and Spinoccia (2018),
 71 see Figure 1A. The exact degraded version of the data is also available (see Data and Code Availability).
 72 Bathymetry is a well-known delineator of marine biodiversity, and so inclusion probabilities are chosen to
 73 vary with depth. To achieve this, the inclusion probabilities are specified by: 1) defining 4 depth bins (Figure
 74 1B), 2) stipulating that within each bin there is the expectation of the same number of samples, and 3)
 75 specifying the inclusion probabilities within each depth bin so sums within bins are equal across bins.

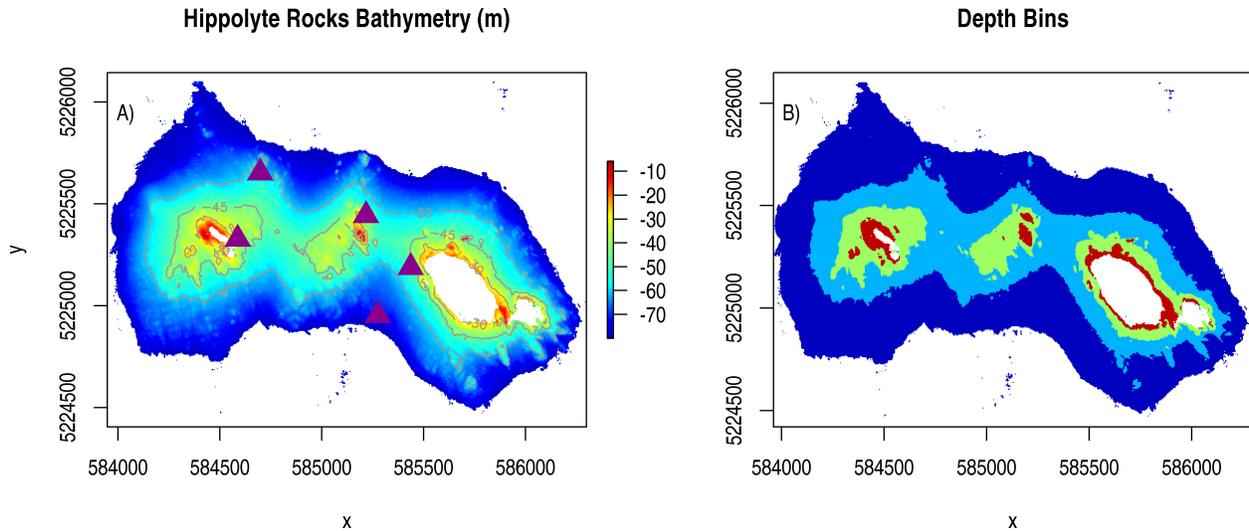


Figure 1: A) Bathymetry of the Hippolyte Rocks survey area. Purple triangles represent the locations of legacy sites that could be incorporated into a survey design. B) The depth bins within which the inclusion probabilities are constant.

76 2.1.1 Generating Point-Based Designs

77 The first design that will be generated is an equal probability spatially-balanced design for nSamp=100
 78 survey sites. This design is appropriate when there is no information about the survey area, including which
 79 locations may have higher variance. With `bathy.df` being a data frame containing the grid of locations
 80 defining the survey area, the call is below, the locations are illustrated in Figure 2A and the first six rows of
 81 the return value is in Table 2.

```
#### A spatially-balanced sample within the study area (not depth-related)
evenSample <- quasiSamp( n=nSamp, potential.sites=bathy.df[,c("x", "y")] )
```

82 The second design is with unequal inclusion probabilities. To generate this design, the inclusion probabilities
 83 are passed to `quasiSamp()`. The resulting design is plotted in Figure 2B and is the result of the call:

Table 2: Four survey sites selected by the function `quasiSamp` for the even inclusion probability design. In order, the columns are: 'x' coordinate (e.g. longitude), 'y' coordinate (e.g. latitude), the inclusion probability that the site was selected with, and the row number from the `potential.sites` input argument. Inclusion probabilities, over all potential sites, will sum to 1 by construction.

x	y	inclusion.probabilities	ID
585996.7	5225038	9.1e-06	78146
584027.6	5225018	9.1e-06	79783
585791.7	5225523	9.1e-06	25518
585590.4	5225588	9.1e-06	18839

```
#### A spatially-balanced sample with shallow locations preferred
unevenSample <- quasiSamp( n=nSamp, potential.sites=bathy.df[,c("x","y")],
                          inclusion.probs=bathy.df[, "inclusion.prob"])
```

84 The third design additionally incorporates the spatial locations of existing legacy sites (Foster *et al.* 2020).
 85 This involves a two-step procedure: 1) alter the inclusion probabilities to down-weight locations near the
 86 legacy sites, and 2) take a sample with the altered inclusion probabilities.

```
#adjust the inclusion probabilities for the locations of the legacy sites
bathy.df$altered.inclusion.prob <- alterInclProbs( as.matrix( legacySites),
          potential.sites=bathy.df[,c("x","y")],
          inclusion.probs=(nSamp-nrow(legacySites))*bathy.df[, "inclusion.prob"],
          mc.cores=8)

#a spatially-balanced sample with legacy sites and a preference for shallow locations
unevenLegacySample <- quasiSamp( n=nSamp-nrow(legacySites),
          potential.sites=bathy.df[,c("x","y")],
          inclusion.probs=bathy.df[, "altered.inclusion.prob"])
```

87 2.1.2 Generating Transect-Based Designs

88 Generating transect-based designs requires special methods as the task of respecting (point-based) inclusion
 89 probabilities using a transect sample is non-trivial. `MBHdesign` employs the method of Foster *et al.* (2020),
 90 who choose transects based on inclusion probabilities of the cells that the transects intersect. Transects are then
 91 chosen based on these derived inclusion probabilities. Complex survey areas and spatially-complex inclusion
 92 probabilities maps will mean that finding transect designs is a difficult and computationally demanding task.
 93 A design for `nSamp=12` transects is shown in Figure 3A, and is generated using:

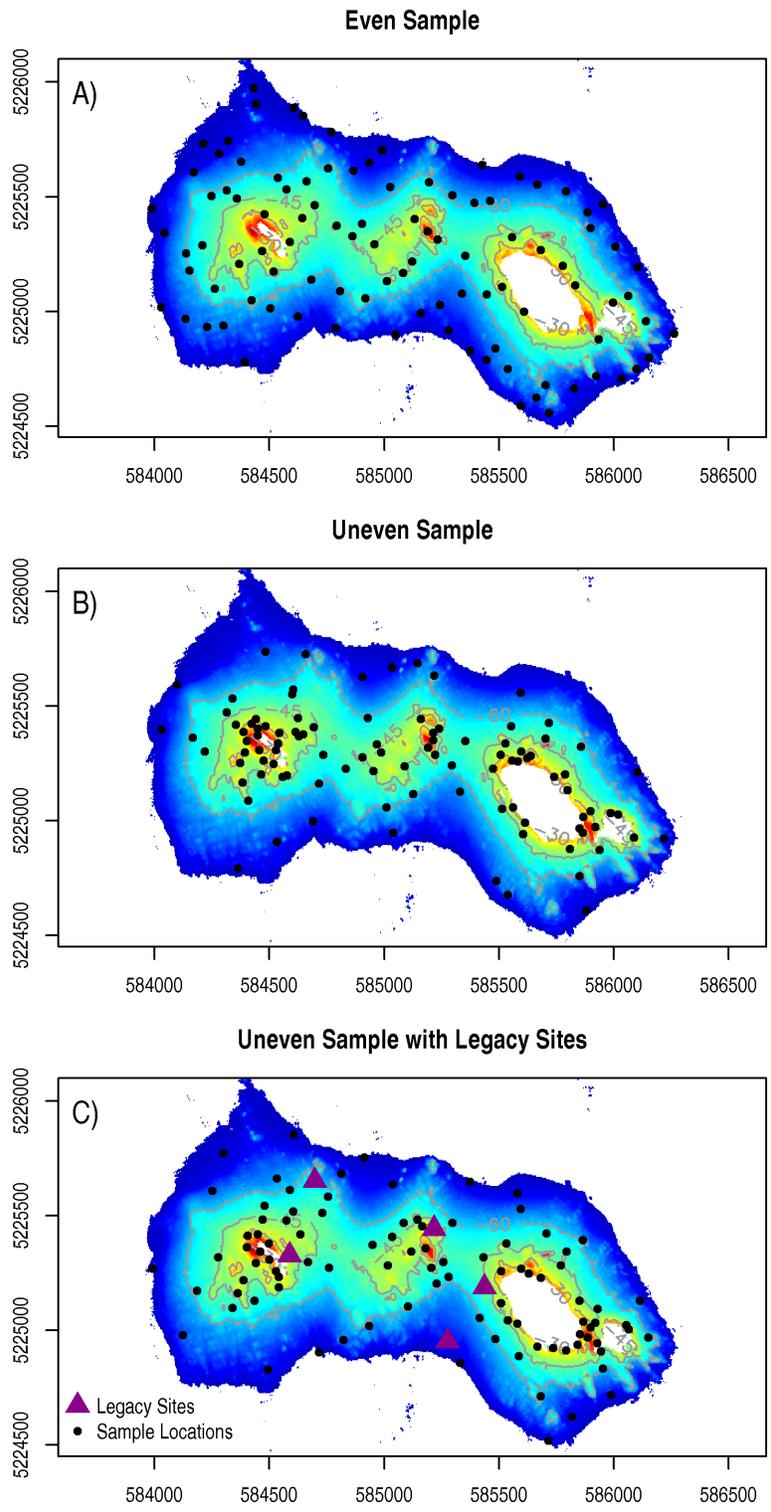


Figure 2: Example point-based designs. A) a spatially-balanced design with equal inclusion probabilities. B) a spatially-balanced design with unequal inclusion probabilities. C) a spatially-balanced design with unequal inclusion probabilities that incorporates 5 legacy sites.

Table 3: Four points (of 120) on transects selected by `transSamp`. Stored in the second element of the return object. In order, the columns are: transect number, the coordinates of the transects' midpoints, the compass bearing of the transect, the points on the transect, the user-defined inclusion probability and the (internal) probability of selection to maintain the user-defined probabilities.

transect	mid_x	mid_y	direction	x	y	specifiedInclProb	AdjustedInclProb
12	584573.0	5225438	51.42857	584579.5	5225443	0.0005543	0.0004217
12	584573.0	5225438	51.42857	584553.4	5225423	0.0006148	0.0000000
11	584298.5	5225288	85.71429	584273.5	5225286	0.0004108	0.0001009
7	585634.4	5224713	51.42857	585614.8	5224698	0.0001113	0.0000662

```
#The representation of transects and other algorithm controls
control <- list( transect.pattern='line', transect.nPts=10, nRotate=21,
                line.length=150, mc.cores=8)

#a spatially-balanced sample with shallow locations preferred
transSample <- transectSamp( n=nSamp, potential.sites=incl.prob[,c("x","y")],
                            inclusion.probs=incl.prob[, "inclusion.prob"], control=control)
```

104 The control argument contains information about the type and representation of the transects, as well as
 105 information about the underlying algorithm. In this example, the 150 m linear transects are represented by
 106 `transect.nPts=10` discrete points and `nRotate=21` different rotations (directions) are considered from each
 107 cell within the survey area. Increasing `transect.nPts` and `nRotate` will increase the ability of the algorithm
 108 to respect the specified inclusion probabilities, but at the expense of computation. The value of `line` for
 109 `transect.pattern` is a special case, in general its value should be a matrix that represents the (centered) shape
 110 of the transect. For a non-linear example, see the supporting information file ‘VisualIllustration6.R’ of Foster
 111 *et al.* (2020).

102 The remaining arguments are similar to those in `quasiSamp()`. The return value of `transectSamp` is a list of
 103 two elements, each of which is structured similarly to the return object of `quasiSamp` (Table 3). The first
 104 element of the `transectSamp()` return list are the `nSamp` central locations of the selected transects. The
 105 columns are the same as in Table 2 but contain additional information about the transect midpoints: the ID
 106 number and the direction from that point. The second element of the return list is for the points on each of
 107 the `nSamp` transects, see Table 3. The two types of inclusion probabilities returned are: those specified by
 108 the user, and those which the algorithm has used to try and obtain transect samples with the user-supplied
 109 probabilities. See Foster *et al.* (2020) for details on the altering process, and see the package vignette for an
 110 illustration of its effect.

111 It may be important to know that a transect’s inclusion probability is just the sum of the inclusion probabilities

112 for the points that represent it. This means that a transect's inclusion probability may be high even though
113 one, or more, of the inclusion probabilities for its constituent points is very low or even zero. This raises
114 a design decision: whether areas just outside the survey area should be excluded completely (inclusion
115 probability of NA) or should they be included if the neighbouring sites are sufficiently important (inclusion
116 probability of 0).

117 It is possible to place constraints on the transects, to remove the chance of obtaining a transect that cannot
118 be executed in the field. As an example, when sampling with towed underwater video (ToV), it is standard
119 to only perform transects that are downhill from the start location. Another example is sampling seamounts,
120 again using ToV, where transects are sometimes taken from a seamount's peak (e.g. Williams *et al.* 2020).
121 There are specific functions within `MBHdesign` to perform both these specific tasks (`findDescendingTrans()`
122 and `findTransFromPoint()`), which return a matrix of text values that describe the behaviour of the transect,
123 which can be then used within `transectSamp()`. Note that with more constraints the inclusion probabilities
124 become less and less well-respected / approximated. Example designs that incorporate these two constraints
125 are given in Figures 3B and 3C.

126 3 Summary and Discussion

127 An ubiquitous and fundamental aspect of field ecology is survey design. To aid generation of efficient designs,
128 the `MBHdesign` R-package is introduced. The designs generated are spatially-balanced and allow for unequal
129 inclusion probabilities. The package offers a simple interface whilst remaining flexible enough to provide a
130 rich suite of designs.

131 Whilst *design* is the primary consideration of `MBHdesign`, *inference* is only briefly considered (using the
132 supplied function `modEsti` for inference with legacy sites). It is noted however, that there are many different,
133 valid, design-based and model-based inferential methods for the data that point-based designs from `MBHdesign`
134 will generate. All approaches will benefit from an efficient design. Foster *et al.* (2020) suggested that model-
135 based analyses for transect-based surveys were necessary, because design-based analyses were currently
136 unknown for this type of randomisation. Such analyses should take into account spatial autocorrelation.

137 The implemented algorithm for generating point-based designs (BAS; Robertson *et al.* 2013) is computationally
138 thrifty. However, the algorithm for transect-based designs (Foster *et al.* 2020) is computationally more
139 expensive. This is partly due to the relative infancy of transect-based approaches, especially compared to
140 points.

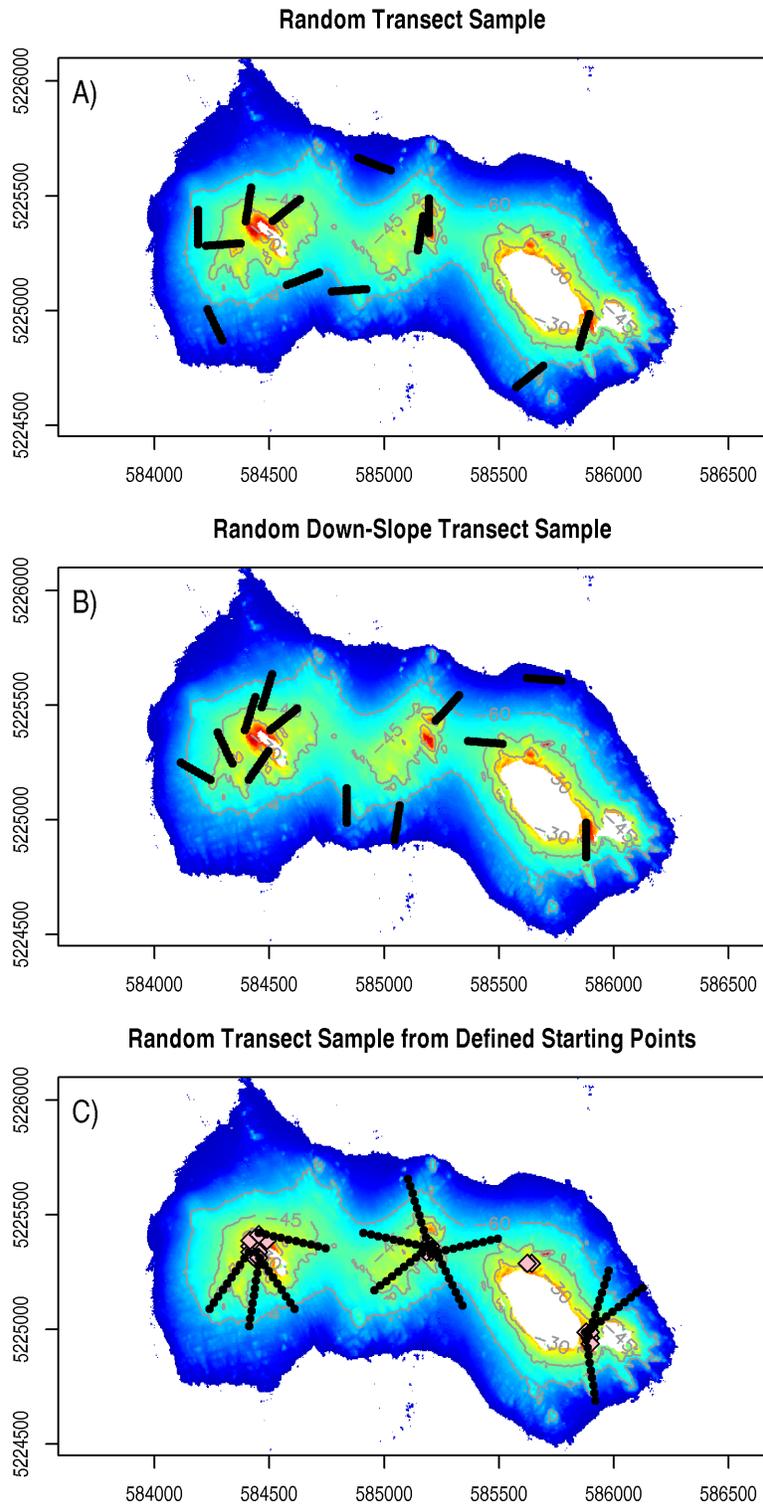


Figure 3: Example transect-based designs. A) A spatially-balanced design with unequal inclusion probabilities. B) Like A) but with transects that are downhill. C) Like A) but only with transects that start in the shallowest 20 grid cells. In all panels, black dots represent way-points on the transect and in C) pink diamonds represent the set of potential start locations.

141 The code-base of `MBHdesign` is freely available from CRAN (<https://cran.r-project.org/package=MBHdesign>)
142 and this article was created using version 2.1.8.

143 **4 Acknowledgements**

144 I would like to thank Emma Lawrence, Jac Monk, Geoff Hosack, Keith Hayes and Skip Woolley for advice and
145 support. This work was undertaken for the Marine Biodiversity Hub, a collaborative partnership supported
146 through funding from the Australian Government’s National Environmental Science Program.

147 **5 Data and Code Availability**

148 The Zenodo repository (Foster 2020) contains the Hippolyte Rocks bathymetry data, and the R-code used to
149 create the examples in this work. The data are a degraded version of that obtained from Spinoccia (2018).
150 The `MBHdesign` package is available from <https://cran.r-project.org/package=MBHdesign>.

151 **6 Author Contributions**

152 S Foster is responsible for all parts of this work.

153 **References**

- 154 Altman, D. (1991). Randomisation. *BMJ*, **302**, 1481.
- 155 Foster, S.D. (2020). Code and data required to run examples in `MBHdesign` description paper. Retrieved
156 from <https://doi.org/10.5281/zenodo.4291227>
- 157 Foster, S.D., Hosack, G.R., Lawrence, E., Przeslawski, R., Hedge, P., Caley, M.J., Barrett, N.S., Williams, A.,
158 Li, J., Lynch, T., Dambacher, J.M., Sweatman, H.P. & Hayes, K.R. (2017). Spatially balanced designs that
159 incorporate legacy sites. *Methods in Ecology and Evolution*, **8**, 1433–1442.
- 160 Foster, S.D., Hosack, G.R., Monk, J., Lawrence, E., Barrett, N.S., Williams, A. & Przeslawski, R. (2020).
161 Spatially balanced designs for transect-based surveys. *Methods in Ecology and Evolution*, **11**, 95–105.
- 162 Grafström, A. (2012). Spatially correlated poisson sampling. *Journal of Statistical Planning and Inference*,

163 **142**, 139–147.

164 Grafström, A. & Lisic, J. (2019). *BalancedSampling: Balanced and spatially balanced sampling*.

165 Grafström, A. & Lundström, N.L.P. (2013). Why well spread probability samples are balanced. *Open Journal*
166 *of Statistics*, **3**, 36–41.

167 Grafström, A., Lundström, N.L.P. & Schelin, L. (2012). Spatially balanced sampling through the pivotal
168 method. *Biometrics*, **68**, 514–520.

169 Hayes, K.R., Hosack, G.R., Lawrence, E., Hedge, P., Barrett, N.S., Przeslawski, R., Caley, M.J. & Foster,
170 S.D. (2019). Designing monitoring programs for marine protected areas within an evidence based decision
171 making paradigm. *Frontiers in Marine Science*, **6**, 746.

172 Hijmans, R.J. (2018). *Raster: Geographic data analysis and modeling*.

173 Kermorvant, C., D’Amico, F., Bru, N., Caill-Milly, N. & Robertson, B. (2019). Spatially balanced sampling
174 designs for environmental surveys. *Environmental Monitoring and Assessment*, **191**.

175 Kincaid, T.M., Olsen, A.R. & Weber, M.H. (2019). *Spsurvey: Spatial survey design and analysis*.

176 Leek, J.T. & Peng, R.D. (2015). What is the question? *Science*.

177 McDonald, T. & McDonald, A. (2020). *SDraw: Spatially balanced samples of spatial objects*.

178 Nichol, S., Anderson, T., McArthur, M., Heap, A., Siwabessy, P. & Brooke, B. (2009). *Southeast Tasmania*
179 *temperate reef survey. Post survey report record 2009/43*. Geoscience Australia, Canberra.

180 Robertson, B.L., Brown, J.A., McDonald, T. & Jaksons, P. (2013). BAS: Balanced acceptance sampling of
181 natural resources. *Biometrics*, **69**, 776–784.

182 Spinoccia, M. (2018). AusSeabed bathymetry holdings. Retrieved July 14, 2020, from [http://pid.geoscience.](http://pid.geoscience.gov.au/dataset/ga/116321)
183 [gov.au/dataset/ga/116321](http://pid.geoscience.gov.au/dataset/ga/116321)

184 Stevens, D. & Olsen, A. (2004). Spatially balanced sampling of natural resources. *Journal of the American*
185 *Statistical Association*, **99**, 262–278.

186 Thompson, S. (2012). *Sampling*. Wiley.

187 Williams, A., Althaus, F., Green, M., Maguire, K., Untiedt, C., Mortimer, N., Jackett, C.J., Clark, M., Bax,
188 N., Pitcher, R. & Schlacher, T. (2020). True size matters for conservation: A robust method to determine
189 the size of deep-sea coral reefs shows they are typically small on seamounts in the southwest pacific ocean.
190 *Frontiers in Marine Science*, **7**, 187.