

# Effects of ignoring survey design information for data reuse

SCOTT D. FOSTER <sup>1,7</sup>, JARNO VANHATALO,<sup>2,3</sup> VERENA M. TRENKEL,<sup>4</sup> TORSTI SCHULZ,<sup>3</sup> EMMA LAWRENCE,<sup>5</sup>  
 RACHEL PRZESLAWSKI,<sup>6</sup> AND GEOFFREY R. HOSACK<sup>1</sup>

<sup>1</sup>Data61 CSIRO, GPO Box 1538, Hobart, TAS 7001 Australia

<sup>2</sup>Department of Mathematics and Statistics, University of Helsinki, P.O. Box 68, Helsinki FIN-00014 Finland

<sup>3</sup>Department of Organismal and Evolutionary Biology Research Program, University of Helsinki, P.O. Box 68, Helsinki FIN-00014 Finland

<sup>4</sup>IFREMER, Rue de l'île d'Yeu, BP 21105, 44311 Nantes Cedex 3, France

<sup>5</sup>Data61 CSIRO, 41 Boggo Rd, Dutton Park, QLD 4102 Australia

<sup>6</sup>Geoscience Australia, GPO Box 378, Canberra, ACT 2601 Australia

*Citation:* Foster, S. D., J. Vanhatalo, V. M. Trenkel, T. Schulz, E. Lawrence, R. Przeslawski, and G. R. Hosack. 2021. Effects of ignoring survey design information for data reuse. *Ecological Applications* 31(6):e02360. 10.1002/eap.2360

**Abstract.** Data are currently being used, and reused, in ecological research at an unprecedented rate. To ensure appropriate reuse however, we need to ask the question: “Are aggregated databases currently providing the right information to enable effective and unbiased reuse?” We investigate this question, with a focus on designs that purposefully favor the selection of sampling locations (upweighting the probability of selection of some locations). These designs are common and examples are those designs that have uneven inclusion probabilities or are stratified. We perform a simulation experiment by creating data sets with progressively more uneven inclusion probabilities and examine the resulting estimates of the average number of individuals per unit area (density). The effect of ignoring the survey design can be profound, with biases of up to 250% in density estimates when naive analytical methods are used. This density estimation bias is not reduced by adding more data. Fortunately, the estimation bias can be mitigated by using an appropriate estimator or an appropriate model that incorporates the design information. These are only available however, when essential information about the survey design is available: the sample location selection process (e.g., inclusion probabilities), and/or covariates used in their specification. The results suggest that such information must be stored and served with the data to support meaningful inference and data reuse.

**Key words:** bias; data; database; findable; accessible; interoperable; reusable data; Horvitz-Thompson estimator; inclusion probability; model; population density estimate; reuse; survey design.

## INTRODUCTION

Ecology and other environmental sciences, like most scientific disciplines, are currently utilizing an unprecedented volume of data (e.g., LaDeau et al. 2017) and are poised to make use of even more (e.g., Culina et al. 2018). In our opinion, this trend is due to two parts: the increase in publicly available databases, and the realization that incorporating data from many sources increases the information available for any particular study (Fletcher et al. 2019). The intended and desirable outcomes from this trend are that individual ecological studies are now broadening their ecological scale (e.g., global studies: Phillips et al. 2019, Gagné et al. 2020, McKenzie et al. 2020), or are shedding brighter lights on smaller scales so that data-poor systems can be quantitatively studied (e.g., Kindsvater et al. 2018, Fletcher et al. 2019).

The quality of the inferences from these analyses is only as good as the data that goes into them (e.g., Dobson et al. 2020). For aggregated data, this means the quality of the contributing data sets and how well they can relate to each other. This is well recognized, and endeavors have been undertaken to improve data quality, with primary focus on two aspects: FAIR (Findable, Accessible, Interoperable, Reusable; Wilkinson et al. 2016, Stall et al. 2019), and standardization of collection methods (e.g., Przeslawski et al. 2019). Undoubtedly, these will increase data reusability. However, are there any other hitherto overlooked aspects that will impede the reusability of ecological data?

All ecological data are the result of some sort of sampling process, and this process is based on a survey plan that describes where and how to collect samples. Many surveys do not consider these aspects in sufficient detail before implementation (Legg and Nagy 2006). Recent modeling efforts with data aggregated from multiple surveys have suggested that survey information, such as the survey plan and sampling gear, should be taken into account to help data “speak” to one another (Fletcher

Manuscript received 29 July 2020; revised 5 November 2020; accepted 4 February 2021. Corresponding Editor: Eric J. Ward.

<sup>7</sup>E-mail: scott.foster@data61.csiro.au

et al. 2019). Without this information, it is hard to understand the meaning of the data and further (potentially wrong) assumptions are required for analysis and interpretation. Indeed, the survey information, or survey metadata, is sometimes not even available to users as the data themselves are. The importance of this omission may be under-appreciated, and it is yet unknown how much of an effect this has on subsequent analyses.

In this work, we investigate what effect ignoring survey design information can have on analysis outputs. We make our inference from a simulation experiment based on a 2018 survey of deep-water corals, which was formally and purposefully designed to increase information content by modifying the selection process for sample locations (Foster et al. 2020). The specific questions we ask are (1) If these data were contributed to databases that aggregate multiple surveys, would naive reuse generate a false picture of the ecology or provide misleading information for management? and (2) How much, if any, modification of the sample location selection process (away from complete randomization) is tolerable before data reuse needs to incorporate survey design information? We discuss what survey design information is needed to be stored within aggregated databases.

## METHODS

### Deep-water corals

A population of the deep-water stony coral *Solenostrea variabilis* is located in the Huon Australia Marine Park, which contains geomorphological features known as the Tasmanian seamounts, located south of Tasmania, Australia. The distribution of *S. variabilis* in this region is not well understood, except in vague terms; it prefers outcropping locations within a partially known depth ranges (Thresher et al. 2011). To rectify this knowledge gap, a scientific survey was undertaken in late 2018 (Williams et al. 2018, 2020), which follows a 2010 survey in a comparable region (Williams et al. 2010). The design for the 2018 survey is outlined in Foster et al. (2020) and consisted of favoring sample locations where *S. variabilis* presence/abundance is thought to be uncertain.

The method used to create the survey was to sample potential sampling locations with specified uneven inclusion probabilities (e.g., Thompson 2012). For the 2018 seamount survey, these probabilities were expert derived and up-weight the locations that (1) are within the broad species bathymetric range and (2) are locally elevated in relation to neighboring locations, measured by the topographic position index (TPI; Weiss 2001); see Fig. 1. Only those locations within 485 and 2,015 m deep were considered for sampling.

In this work, we utilize the 2018 survey's uneven inclusion probabilities defined in Foster et al. (2020, Table 2), which links our simulation to procedures used in practice. These inclusion probabilities are highly skewed as the area covered by seamounts is comparatively small

(See Fig. 1). The distribution of inclusion probabilities is given in Appendix S1: Fig. S2. To simplify computation, we only use the survey area within the Huon Marine Park, which also contains many of the seamounts in the broader region.

We also utilize data on *S. variabilis* from a 2010 survey described in Williams et al. (2010). The survey design for the 2010 survey was less formal but did target the coral's depth range and sites with higher TPI. For modeling purposes, we assume that the 2010 design is *ignorable* once the depth and TPI are included as covariates (Gelman et al. 2013). These data were generated from a camera towed along the seafloor, and later quantified by counting the number of live *S. variabilis* coral heads within regularly spaced images. The size of the seafloor covered by the quantification area, within each image is also recorded. Overall, in the Huon park there are 1,517 images spaced along 19 transects with the longest transect having 212 images and the shortest 12. Images from the 2018 survey were not used in this work as, at the time of writing, the images are not yet quantified.

### A model for coral distribution

To analyze the 2010 image data, we use a geostatistical model. In particular, we use the "SPDE" approach, which is implemented using the "INLA" approximation (Rue et al. 2009, Lindgren and Rue 2015) implemented for R (R Core Team 2019). This approach to computing is relatively fast, so that many models can be fitted. We notate each of the ( $i = 1 \dots 1,517$ ) observed *S. variabilis* coral abundance data as  $y_i$ , and model all observations as a function of geographical position ( $s_i$ ), bathymetry, and TPI. That is

$$\log[E(y_i | \boldsymbol{\theta}, b(s_i), t(s_i))] = \beta_0 + \beta_1 b(s_i) + \beta_2 b(s_i)^2 + \beta_3 t(s_i) + u(s_i) + \log(A_i), \quad (1)$$

where  $\beta_j$  is a regression parameter,  $b(s_i)$  and  $t(s_i)$  are bathymetry and TPI covariates, respectively,  $u(s_i)$  is a spatial random variable,  $A_i$  is the area that the  $i$ th image sampled, and all effects are gathered into the parameter vector  $\boldsymbol{\theta}$ . A quadratic effect for depth was assumed to reflect the belief that the *S. variabilis* depth-niche was covered by the data, whereas it is thought that there is no upper limit to TPI preference. We assume that the conditional distribution of  $y_i | \boldsymbol{\theta}, b(s_i), t(s_i)$  is Poisson and that the spatial random variable,  $u(s_i)$ , is assumed to follow a Matérn Gaussian process with mean zero and smoothness  $\nu = 1$ . This model gives the spatial covariance of the random effect as

$$\text{cov}[u(s_i), u(s_{i'})] = \sigma^2 \kappa \left( \kappa |s_i - s_{i'}| \right) K_1 \left( \kappa |s_i - s_{i'}| \right),$$

which has standard deviation ( $\sigma$ ) and scaling parameter ( $\kappa$ ). The function  $K_1(\cdot)$  is the modified Bessel function of the second kind and order 1. The Matérn process has

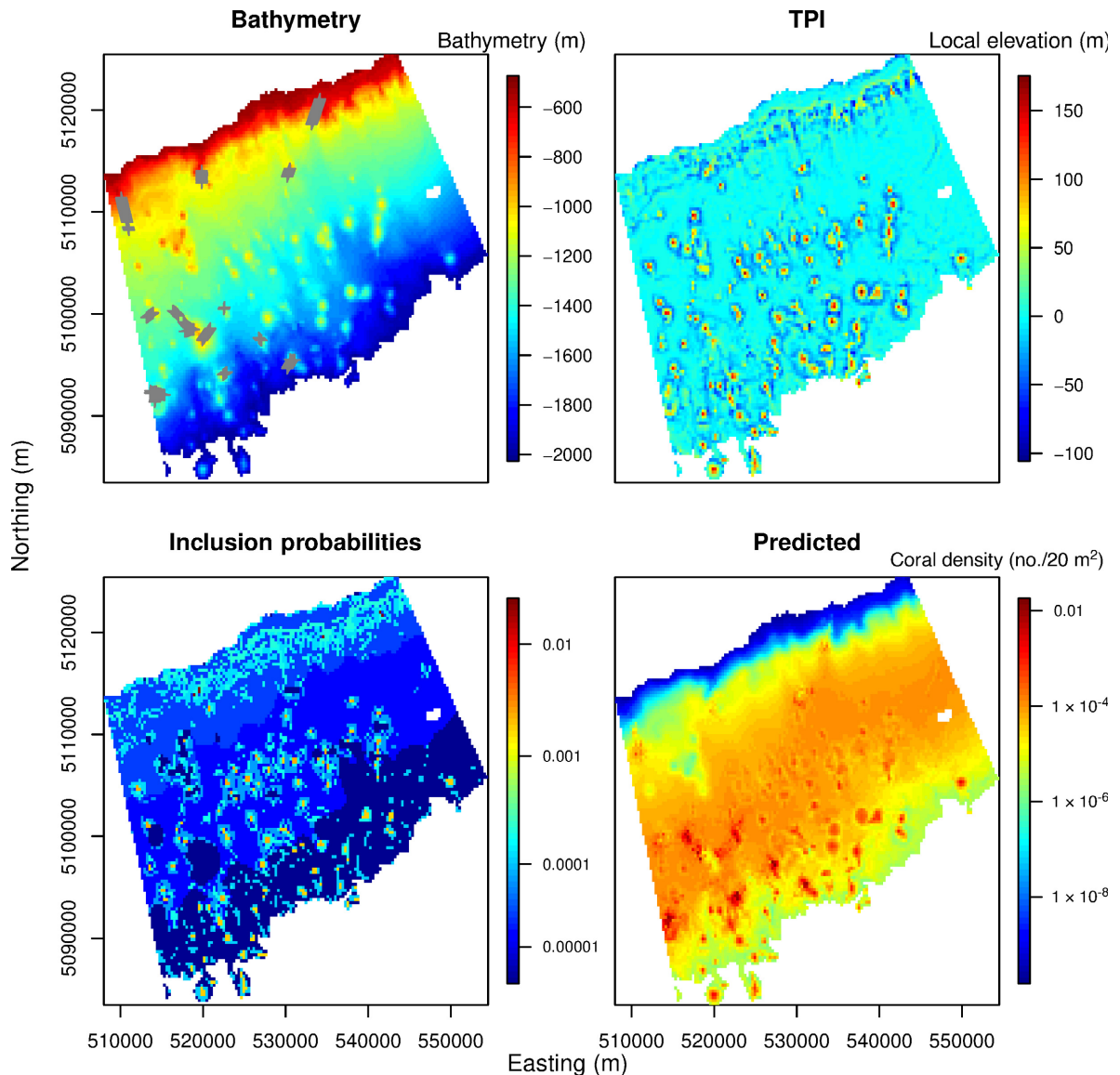


FIG. 1. Detail of the sampling locations within the Huon Australian Marine Park, located south of Tasmania, Australia. These locations are those that are within the depth range of 485 and 2,015 m. Bathymetry is water depth (m), and TPI is topographic position index and gives an indication of how elevated each cell is with respect to its neighbors (units of TPI are meters). The inclusion probabilities are those used to draw the sampling locations for the survey. The predicted values are from the model defined in *A Model for Coral Distribution*, fitted to the original survey data whose locations are gray “+” on the bathymetry map. The image-frame size for prediction ( $20 \text{ m}^2$ ) is arbitrary. The coordinate reference system used is WGS 84/UTM zone 55S, with units of m east and north.

effective range of  $\sqrt{8}/\kappa$ , which is the empirically derived spatial distance where correlation is  $\gamma \approx 0.1$  (Lindgren et al. 2011, Lindgren and Rue 2015). We specify a penalized complexity prior (Simpson et al. 2017) where there is  $\Pr(\sigma > 5) = 0.1$ , which penalizes overly flexible spatial processes. The effective range ( $\gamma$ ) of the process has a prior such that  $\Pr(\gamma < 50 \text{ m}) = 0.05$  so that the spatial dependence is unlikely to be very short. Priors for the regression coefficients are chosen to penalize extreme values. We define these to be normal distributions with

zero mean and variance equal to 5. Both covariates were standardized to have mean zero and variance 1 before analysis.

#### Simulation experiment

The base form of the simulation experiment is (1) vary inclusion probabilities to be more and less severe than the 2018 inclusion probabilities, (2) generate a survey design from these inclusion probabilities, (3) simulate

data at the sampling locations generated (using the model fitted to the 2010 image data), (4) analyze the simulated data with naive (ignoring sampling probabilities) and more sophisticated methods that account for the survey design, and (5) summarize the simulations' analyses as a response to variation in the unevenness of inclusion probabilities. This approach will inform if the survey data can be naively reused in the analysis of aggregated data.

To vary the inclusion probabilities for the  $N = 8,840$  sites that define the sampling area, we start with the inclusion probabilities used to design the 2018 survey, and we arrange these probabilities into an  $N \times 1$  vector  $\mathbf{p}$ . The  $N$  sites are arranged on a  $300 \times 300$  m grid and match the grid of the covariates (see Fig. 1). This was chosen to match that used in Foster et al. (2020), who used this as a compromise between accuracy and computational expense. The inclusion probabilities for the as

$$\mathbf{p}_\alpha = \max(\mathbf{p}_\alpha^*, 0)/K,$$

where

$$\mathbf{p}_\alpha^* \triangleq [\mathbf{1}\bar{p} + \alpha(\mathbf{p} - \mathbf{1}\bar{p})],$$

$\bar{p} = (\mathbf{1}^T \mathbf{p})/N$  is the mean of  $\mathbf{p}$ ,  $K = \mathbf{1}^T \mathbf{p}_\alpha$  is a normalizing constant, and the maximum function is applied element-wise. If an inclusion probability is zero, then that site will not be chosen in the sample. The parameter  $\alpha$  indexes the severity of the unevenness in the inclusion probabilities, with  $\alpha = 0$  corresponding to even inclusion probabilities (and completely randomized sampling),  $\alpha = 1$  corresponding to the 2018 survey's inclusion probabilities and  $\alpha > 1$  giving inclusion probabilities more extreme. We allow  $\alpha$  to vary from 0 to 2 in increments of 0.1. For each  $\alpha$ ,  $J = 1,000$  surveys were simulated, each consisting of  $n = 50,100,200$  observations from the  $N$  sites within the sampling area. The locations of the observations were chosen at random using  $\mathbf{p}_\alpha$ .

For each simulated survey, data were simulated at the  $n$  selected locations using parameters drawn from the posterior distribution of the model in *A Model For Coral Distribution*, fitted to the 2010 data. This ensures that all modeled aspects of the 2010 data, including variability, are incorporated into the simulation study. The marginal posterior distribution of the covariate effects is presented in Appendix S1: Fig. S3.

Each simulated data set is analysed using design-based and model-based estimators. The target metric in each of these analyses is the average number of corals per  $20\text{-m}^2$  image (coral density). Theoretically, it is useful to consider the bias in the average density for both design-based and model-based analyses: design-based estimates are intended to be unbiased for the average, and the average is also the Bayes estimate under quadratic loss for model-based methods. We note that other summaries could be of interest, like the maximum coral

density, but the average is a very common summary, almost ubiquitously so. The design-based analyses were a naive mean ( $\ln \sum y_i$ ), and the Horvitz-Thompson (HT) estimator (see Thompson 2012) of the form  $\sum y_i n p_{\alpha i}$  where the sum is over the  $n$  samples. The HT estimator is only available when the inclusion probabilities for the samples are known, and it should (theoretically) produce unbiased estimates, even when inclusion probabilities are unequal. The naive mean should (theoretically) only be unbiased when the inclusion probabilities are equal (Thompson 2012).

The model in *A Model for Coral Distribution* was used to analyze each simulated data set along with three simplifications. These models are used to investigate the effect of only making part of the design information available to the analysis process. The models are as follows:

Covariates + Spatial. The full model in *A Model for Coral Distribution*.

Spatial. Covariates unavailable or neglected and only the spatial effects are included.

Covariates. Spatial effects are omitted. The analyst assumes that the observations are independent given the covariates.

Bathymetry/TPI. The third simplification is to drop each of the covariates (bathymetry and TPI) in turn, with no spatial effect.

For all models, the "true" average density of the  $j$ th simulation,  $\mu_j$ , was calculated by taking the mean of the set of predictions formed at a grid of  $N$  locations throughout the study region. The same set of draws of the parameters (from the posterior that conditions on the 2010 data, *A Model for Coral Distribution*) were used to calculate the set of  $\mu_j$ . For a given value of  $\alpha$ , the average density estimate of the  $k$ th estimation method was assessed by calculating a percentage difference between the estimated average density ( $\hat{\mu}_{jk}$ ) and the quantity it is estimating ( $\mu_j$ ). Formally, for the  $j$ th simulation replicate and the  $k$ th estimation method, the percentage difference is

$$d_p(j, k) = 100 \frac{\hat{\mu}_{jk} - \mu_j}{\mu_j}.$$

For each value of  $\alpha$  and for each estimation method, there are  $J$  estimates of average coral density. We summarize this information using the median and mean absolute deviation [MAD; Venables and Ripley 2002]. These are relatively robust measures of location and scale that are not unduly affected by extreme values (outliers). We take the median of the naive mean estimates, when the inclusion probabilities were even ( $\alpha = 0$ ), as the reference value for comparison against all other estimators and all other values of  $\alpha$ . The naive mean has well known and desirable properties when sampling is even ( $\alpha = 0$ ).

RESULTS

Fitting the model to the 2010 image data, see *A Model for Coral Distribution*, suggested that coral density peaked around 1,350 m deep, and had a much reduced expectation outside of the range (-1,700 to -1,000 m). Increasing TPI increased the density of corals (about 12 times increase from flat areas to the extremely elevated). The spatial dependence was short with  $E(\gamma|y) = 333$  m ( $SD(\gamma|y) = 72.3$  m), and the spatial standard deviation

was  $E(\sigma|y) = 2.8$  ( $SD(\sigma|y) = 0.4$ ). Posterior distributions for all parameters defined in (1) are presented in Appendix S1: Fig. S3. Posterior predictions from this model are presented in Fig. 1 and show the effect of depth, which is smooth over the survey area, and the relatively patchy effects of TPI and spatial noise.

Results for the simulation experiment, described in *Simulation experiment*, are presented in Fig. 2. Overall, it is clear that ignoring the inclusion probability information can induce substantial bias in average coral

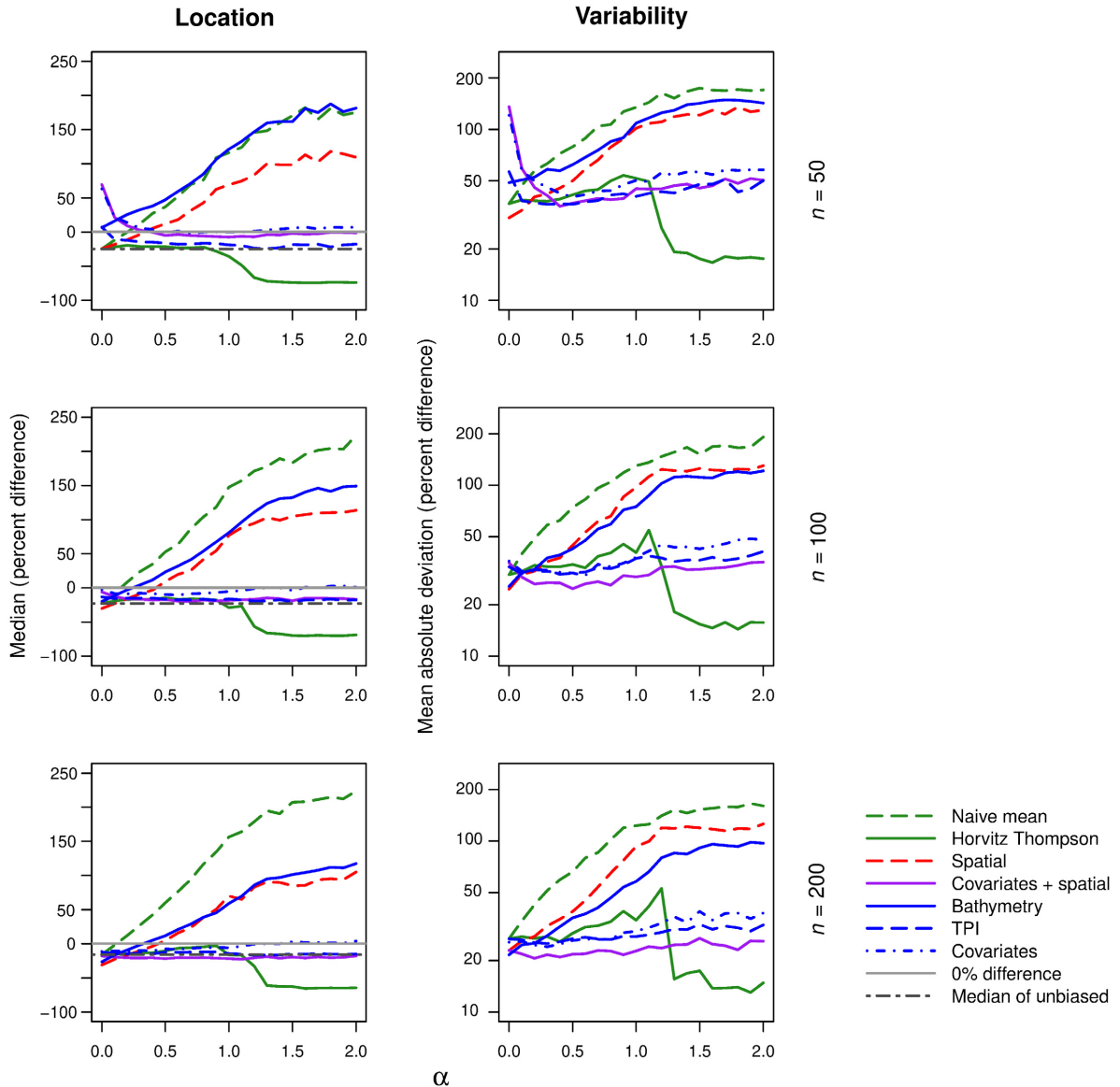


FIG. 2. Results of the simulation experiment based on the survey of the Huon Australian Marine Park. Top row is for surveys with  $n = 50$  sample locations, middle row with  $n = 100$ , and bottom row with  $n = 200$ . Left panels give, for each method and for each  $\alpha$ , the median of the estimates from each of the  $J = 1,000$  simulated data sets. Right panels show the mean absolute deviation (MAD) estimate of variation of the same estimates. See *Methods* for the definition of percent difference and for the choice of reference. Solid gray line is 0% difference and dashed gray line is the median of the naive mean at  $\alpha = 0$  (an unbiased estimator). Small values of  $\alpha$  give more even inclusion probabilities.

density estimates. It is evident though, that even those estimation methods that do incorporate inclusion probabilities can perform badly but in general they work as intended (Fig. 2).

The naive mean is an increasing function of  $\alpha$ , implying that the mean increases as more favorable environments are sampled with greater inclusion probabilities. The naive mean also has very high variation, presumably due to not taking the appropriate weighting of each observation. The HT estimator, which does account for unequal inclusion probabilities, *decreased* with  $\alpha$  and did so sharply just past  $\alpha = 1$  after agreeing with the reference well for all sample sizes for  $\alpha < 1$ .

The simulation illustrated that model-based analyses can produce unbiased estimates of the average density (Fig. 2). The form of the model appears to be important though. The model with no covariates (just a spatial term) and the model with only the bathymetry covariate had undesirable performance, with a trend similar to, but not as extreme as, the naive mean estimate (Fig. 2). When the full model (covariates and spatial) and the TPI-only model were used to analyze the simulated data sets, the median of the estimates for average density were comparatively unbiased albeit after having high values for very small  $\alpha$  with  $n = 50$  (Fig. 2). A similar pattern was observed for the model with both covariates, but this exhibited a slight positive bias.

The full model (with random spatial effects) consistently exhibits small variation in the distribution of estimates, except for  $n = 50$  and for small  $\alpha$  (Fig. 2, right column). This result is linked to the extrapolation/leverage issues (see *Discussion*). The covariates model and the TPI model also suffer from this behavior, at  $n = 50$  and  $\alpha = 0$ , but do not have the low variability in the distribution of estimates, which is exhibited by the full model.

#### SUMMARY AND DISCUSSION

For data to be FAIR it must be reusable (Wilkinson et al. 2016, Stall et al. 2019). For it to be reusable, the relevant information must be made available about *how* to reuse it. Without this information assumptions must be made, with the naive assumption (equal probability random sample) often being wrong.

In this study, we investigated the effect of ignoring survey-design information using a simulation experiment based on a 2018 survey design, and 2010 image data, for a chain of seamounts in southern Australia. We found that ignoring survey design information can induce a substantial bias in estimates of average population density when a naive or an inappropriate analysis method is used; the median of the simulations' average density estimates can be up to 250% biased and estimates for individual data sets even worse. The potentially large bias has the potential to make seemingly straightforward inferences wrong and misleading. We note that the density bias does not disappear with increased sample sizes

(Fig. 2), so "big-data" are no panacea. Even worse, big-data may lead to confident, but biased, inferences.

The simulation experiment showed that some analysis methods performed better than others with uneven inclusion probabilities. The naive mean estimate for population density was the worst performer and some model-based estimators also produced consistently poor results (Fig. 2). The bias was alleviated by incorporating survey design information into the analysis, either through inclusion probabilities for the Horvitz-Thompson (HT) estimator, or through inclusion of the appropriate covariates in a model-based analysis. The sudden appearance of bias in the HT estimator at  $\alpha = 1$  is suspected to be caused by the introduction of sites with inclusion probabilities of zero at  $\alpha = 1$  (see *Methods*) and the associated severe right skew in the distribution of inclusion probabilities (Appendix S1: Fig. S2). We stress that obtaining bias by ignoring design information is not a new result, see Gelman et al. (2013: Chapter 8), Diggle et al. (2010), and Pati et al. (2011). However, this is perhaps under-appreciated by those who deal with ecological data (but see Pennino et al. 2018, Dobson et al. 2020). In fisheries, the problem is receiving recent attention for commercial catch data (e.g., Trenkel et al. 2013).

The poor performance of the models with covariates for smaller sample sizes is likely to be due to insufficient sampling of covariate space (top panel of Appendix S1: Fig. S1,  $\alpha \lesssim 0.2$ ). The insufficient sampling of covariates potentially leads to survey data that must be extrapolated, in covariate space, to predict to all locations (to calculate the average density). This extrapolation in covariates may be erratic and of low quality. The poor sampling of covariates potentially also leads to samples that have undue leverage, which can distort the model estimates. The supplementary study in Appendix S1: Section S1 indicates that small sample sizes underestimate the range of both the bathymetry and TPI covariates. A second reason for poor performance is poor sampling of the spatial extent and hence poor prediction of the spatial random effect throughout the entire region. However, the spatial effect has a relatively small effective range so it is likely that only the largest sample sizes will cover the area sufficiently.

Survey designs are often based on covariates. To account for the influence of the survey design on the model's predictions, these covariates should be included in any model utilizing the survey data (Gelman et al. 2013). If there is no information about how the survey was designed, then it may be most appropriate to include the covariates that the analysts *assumes* to be important in the design, or to use a preferential sampling model (Diggle et al. 2010). We stress that not including any covariates makes the assumption that there were no design-covariates, corresponding to the naive mean in our simulation study, which may be a very inappropriate assumption. We are also aware that this simple advice may be hard to implement in certain situations; an

example is when all covariates are not available for all surveys utilized in a particular reuse. In these situations, careful and skillful analyses must be undertaken, which will rest on assumptions that are necessary to describe *both* the sampling process *and* ecological processes (Diggle et al. 2010, Pati et al. 2011, Liu and Vanhatalo 2020). We note that including a spatial random effect in the southern seamount simulation is not an effective replacement for covariates and that all the design covariates need to be included (Fig. 2). Both these results are likely to be due to the relatively noisy, patchy and spatially non-smooth geographical distribution of TPI.

The southern seamount survey example is quite extreme in its patchy topography and hence the unevenness of the inclusion probabilities. This is why we chose this survey design: to investigate how bad things could be if ignored. However, altering the amount of unevenness (varying  $\alpha$ , *Simulation experiment*) and coupling to the more general theoretical results (e.g., Diggle et al. 2010, Gelman et al. 2013) suggest that our results are generalizable to any survey. Of course, the severity will depend on the amount of variation in the inclusion probabilities, the sample size (Fig. 2), and the survey design (through specification of inclusion probabilities/strata, Fig. 2).

To ensure the ability to reuse data, we suggest that database managers should facilitate the storage and serving of information about survey design, perhaps even incorporated into formal data formats. Reusers of data should be encouraged, perhaps by changing default function settings, to download this information with the data. Data reusers should also be educated about the importance of survey design information. To be clear, this information at minimum should consist of a detailed description of, or accurate reference to, the survey design procedure. Additionally, it is highly desirable to also include (1) the inclusion probabilities (the H-T estimator only needs these at the *sampled* locations) and (2) the values of the covariates at each location within the well-defined study region. We note that the inclusion probabilities could be stored as a field in the data (architecturally similar to another biological measurement), and that the covariates could be part of a meta-data record (or a link to them).

A corollary to this work is that it is best, and in many ways practically necessary, to have a formal survey design if the data are to be reused. While it is possible to model the data from surveys without formal designs, the process becomes more complex (see the variety of models in Diggle et al. 2010 and Gelman et al. 2013: Chapter 8), and is liable to ambiguity through the necessity of making assumptions that are oftentimes untestable. The data may end up being unusable, produce ambiguous results, and their curation and analysis may create a large, hidden research cost (Dobson et al. 2020). We recommend that surveys should be formally designed *and importantly*: the survey design should be stored along with the data. This work serves as a cautionary tale for those who wish to use and reuse data: Do not ignore how the data were obtained, unless you are confident

that there is no intentional, or unintentional, specification of unequal inclusion probabilities in the survey design. Further, this work demonstrates what is needed to interpret survey data: information about the survey design employed to collect the data.

#### ACKNOWLEDGMENTS

This work was undertaken for the Marine Biodiversity Hub, a collaborative partnership supported through funding from the Australian Government's National Environmental Science Program. J. Vanhatalo was additionally funded by the Academy of Finland (grant 317255). R. Przeslawski publishes with the permission of the CEO, Geoscience Australia. We would like to thank Franzis Althaus, Alan Williams, Tim Langlois, Zhi Huang, Jasmine Bursic, Nicholas Johannsohn, Amy Nau, Keith Hayes, and two anonymous reviewers.

#### LITERATURE CITED

- Culina, A., T. Crowther, J. Ramakers, P. Gienapp, and M. Visser. 2018. How to do meta-analysis of open datasets. *Nature Ecology and Evolution* 2:1053–1056.
- Diggle, P. J., R. Menezes, and T.-L. Su. 2010. Geostatistical inference under preferential sampling. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 59:191–232.
- Dobson, A., et al. 2020. Making messy data work for conservation. *One Earth* 2:455–465.
- Fletcher, R. J. Jr., T. J. Hefley, E. P. Robertson, B. Zuckerberg, R. A. McCleery, and R. M. Dorazio. 2019. A practical guide for combining data to model species distributions. *Ecology* 100:e02710.
- Foster, S. D., G. R. Hosack, J. Monk, E. Lawrence, N. S. Barrett, A. Williams, and R. Przeslawski. 2020. Spatially balanced designs for transect-based surveys. *Methods in Ecology and Evolution* 11:95–105.
- Gagné, T. O., G. Reygondeau, C. N. Jenkins, J. O. Sexton, S. J. Bograd, E. L. Hazen, and K. S. Van Houtan. 2020. Towards a global understanding of the drivers of marine and terrestrial biodiversity. *PLoS ONE* 15:1–17.
- Gelman, A., J. Carlin, H. Stern, D. Dunson, A. Vehtari, and D. Rubin. 2013. *Bayesian data analysis*. Third edition. Chapman & Hall/CRC Texts in Statistical Science. Taylor & Francis, Boca Raton, Florida, USA.
- Kindsvater, H. K., N. K. Dulvy, C. Horswill, M.-J. Juan-Jordá, M. Mangel, and J. Matthiopoulos. 2018. Overcoming the data crisis in biodiversity conservation. *Trends in Ecology & Evolution* 33:676–688.
- LaDeau, S., B. Han, E. Rosi-Marshall, and K. Weathers. 2017. The next decade of big data in ecosystem science. *Ecosystems* 20:274–283.
- Legg, C. J., and L. Nagy. 2006. Why most conservation monitoring is, but need not be, a waste of time. *Journal of Environmental Management* 78:194–199.
- Lindgren, F., and H. Rue. 2015. Bayesian spatial modelling with  $r$ -inla. *Journal of Statistical Software, Articles* 63:1–25.
- Lindgren, F., H. Rue, and J. Lindström. 2011. An explicit link between Gaussian fields and Gaussian Markov random fields: the stochastic partial differential equation approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 73:423–498.
- Liu, J., and J. Vanhatalo. 2020. Bayesian model based spatiotemporal survey designs and partially observed log Gaussian cox process. *Spatial Statistics* 35:100392.
- McKenzie, L., L. M. Nordlund, B. L. Jones, L. C. Cullen-Unsworth, C. M. Roelfsema, and R. Unsworth. 2020. The

- global distribution of seagrass meadows. *Environmental Research Letters* 15:074041.
- Pati, D., B. J. Reich, and D. B. Dunson. 2011. Bayesian geostatistical modelling with informative sampling locations. *Biometrika* 98:35–48.
- Pennino, M., I. Paradinas, J. Illian, F. Muñoz, J. Bellido, A. López-Quílez, and D. Conesa. 2018. Accounting for preferential sampling in species distribution models. *Ecology and Evolution* 9:653–663.
- Phillips, H. R. P., et al. 2019. Global distribution of earthworm diversity. *Science* 366:480–485.
- Przeslawski, R., S. Foster, J. Monk, N. Barrett, P. Bouchet, A. Carroll, T. Langlois, V. Lucieer, J. Williams, and N. Bax. 2019. A suite of field manuals for marine sampling to monitor Australian waters. *Frontiers in Marine Science* 6:177.
- R Core Team. 2019. R: a language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria.
- Rue, H., S. Martino, and N. Chopin. 2009. Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 71:319–392.
- Simpson, D. P., H. Rue, A. Riebler, T. G. Martins, and S. H. Sørbye. 2017. Penalising model component complexity: A principled, practical approach to constructing priors. *Statistical Science* 32:1–28.
- Stall, S., L. Yarmey, J. Cutcher-Gershenfeld, B. Hanson, K. Lehnert, B. Nosek, M. Parsons, E. Robinson, and W. Lesley. 2019. Make all scientific data fair. *Nature* 570:27–29.
- Thompson, S. 2012. *Sampling*. Wiley, Hoboken, New Jersey, USA.
- Thresher, R. E., J. Adkins, S. J. Fallon, K. Gowlett-Holmes, F. Althaus, and A. Williams. 2011. Extraordinarily high biomass benthic community on southern ocean seamounts. *Scientific Reports* 1:119.
- Trenkel, V. M., J. A. Beecham, J. L. Blanchard, C. T. T. Edwards, and P. Lorance. 2013. Testing cpue-derived spatial occupancy as an indicator for stock abundance: application to deep-sea stocks. *Aquatic Living Resources* 26:319–332.
- Venables, W., and B. Ripley. 2002. *Modern applied statistics with S*. Fourth edition. Springer, Berlin, Germany.
- Weiss, A. 2001. Topographic positions and landforms analysis (poster). ESRI International User Conference July 2001, ESRI, San Diego, California, USA.
- Wilkinson, M. D., et al. 2016. The FAIR Guiding Principles for scientific data management and stewardship. *Scientific Data* 3:160018.
- Williams, A., et al. 2020. True size matters for conservation: A robust method to determine the size of deep-sea coral reefs shows they are typically small on seamounts in the southwest pacific ocean. *Frontiers in Marine Science* 7:187.
- Williams, A., N. Bax, M. Clark, and T. Schlacher. 2018. RV Investigator voyage summary IN2018 v08: Status and recovery of deep-sea coral communities on seamounts in iconic Australian marine reserves. Australian Marine National Facility Report. [https://www.marine.csiro.au/data/reporting/get\\_file.cfm?eov\\_pub\\_id=187](https://www.marine.csiro.au/data/reporting/get_file.cfm?eov_pub_id=187)
- Williams, A., T. A. Schlacher, A. A. Rowden, F. Althaus, M. R. Clark, D. A. Bowden, R. Stewart, N. J. Bax, M. Consalvey, and R. J. Kloser. 2010. Seamount megabenthic assemblages fail to recover from trawling impacts. *Marine Ecology* 31:183–199.

## SUPPORTING INFORMATION

Additional supporting information may be found online at: <http://onlinelibrary.wiley.com/doi/10.1002/eap.2360/full>