

1 **Running title:** Spatial models for distance sampling  
2 **Number of words:** ~5284  
3 **Number of tables:** 0  
4 **Number of figures:** 6  
5 **Number of references:** 47

6 **Spatial models for distance sampling data:**  
7 **recent developments and future directions**

8 **David L. Miller<sup>1\*</sup>, M. Louise Burt<sup>2</sup>,**  
9 **Eric A. Rexstad<sup>2</sup>, Len Thomas<sup>2</sup>.**

- 10 *1. Department of Natural Resources Science, University of Rhode Island,*  
11 *Kingston, Rhode Island 02881, USA*  
12 *2. Centre for Research into Ecological and Environmental Modelling,*  
13 *The Observatory, University of St. Andrews, St. Andrews KY16 9LZ, UK*

14 **\*Correspondence author. [dave@ninepointeightone.net](mailto:dave@ninepointeightone.net)**

## Summary

1. Our understanding of a biological population can be greatly enhanced by modelling their distribution in space and as a function of environmental covariates. Such models can be used to investigate the relationships between distribution and environmental covariates as well as reliably estimate abundances and create maps of animal/plant distribution.
2. Density surface models consist of a spatial model of the abundance of a biological population which has been corrected for uncertain detection via distance sampling methods.
3. We review recent developments in the field and consider the likely directions of future research before focussing on a popular approach based on generalized additive models. In particular, we consider spatial modelling techniques that may be advantageous to applied ecologists such as quantification of uncertainty in a two-stage model and smoothing in areas with complex boundaries.
4. The methods discussed are available in an R package developed by the authors (**dsm**) and are largely implemented in the popular Windows software Distance.

**Keywords:** abundance estimation, Distance software, generalized additive models, line transect sampling, point transect sampling, population density, spatial modelling, wildlife surveys

## 39 Introduction

40 When surveying biological populations it is increasingly common to record  
41 spatially referenced data, for example: coordinates of observations, habitat  
42 type, elevation or (if at sea) bathymetry. Spatial models allow for vast data-  
43 bases of spatially-referenced data (e.g. OBIS-SEAMAP, Halpin *et al.*, 2009)  
44 to be harnessed, enabling investigation of interactions between environmental  
45 covariates and population densities. Mapping the spatial distribution of a  
46 population can be extremely useful, especially when communicating results  
47 to non-experts. Recent advances in both methodology and software have  
48 made spatial modelling readily available to the non-specialist (e.g., Wood,  
49 2006; Rue *et al.*, 2009). Here we use the term “spatial model” to refer to  
50 any model that includes any spatially referenced covariates, not only those  
51 models that include location as a covariate. This article is concerned with  
52 combining spatial modelling techniques with distance sampling (Buckland  
53 *et al.*, 2001, 2004).

54 Distance sampling extends plot sampling to the case where detection  
55 is not certain. Observers move along lines or visit points and record the  
56 distance from the line or point to the object of interest ( $y$ ). These distances  
57 are used to estimate the *detection function*,  $g(y)$  (for example, Fig. 1), by  
58 modelling the decrease in detectability with increasing distance from the  
59 line or point (conventional distance sampling, CDS). The detection function  
60 may also include covariates (multiple covariate distance sampling, MCDS;  
61 Marques *et al.*, 2007) which affect the scale of the detection function. From  
62 the fitted detection function, the average probability of detection can be

63 estimated by integrating out distance. The estimated average probability  
 64 that an animal is detected given that it is in the area covered by the survey,  
 65  $\hat{p}_i$ , can then be used to estimate abundance as

$$\hat{N} = \frac{A}{a} \sum_{i=1}^n \frac{s_i}{\hat{p}_i}, \quad (1)$$

66 where  $A$  is the area of the study region,  $a$  is the area covered by the survey  
 67 (i.e., the sum of the areas of all of the strips/circles) and the summation  
 68 takes place over the  $n$  observed clusters, each of size  $s_i$  (if individuals are  
 69 observed,  $s_i = 1 \forall i$ ) (Buckland *et al.*, 2001, Chapter 3). Often up to half  
 70 the observations in a plot sampling data set are discarded to ensure the  
 71 assumption of certain detection is met. In contrast, distance sampling uses  
 72 observations that would have been discarded to model detection (although  
 73 typically some detections are discarded beyond a given *truncation distance*  
 74 during analysis).

75 Estimators such as eqn (1) rely on the design of the study to ensure  
 76 that abundance estimates over the whole study area (scaling up from the  
 77 covered region) are valid. This article focusses on *model-based* inference  
 78 to extrapolate to a larger study area. Specifically, we consider the use of  
 79 spatially explicit models to investigate the response of biological populations  
 80 to biotic and abiotic covariates that vary over the study region. A spatially-  
 81 explicit model can explain the between-transect variation (which is often a  
 82 large component of the variance in design-based estimates) and so using a  
 83 model-based approach can lead to smaller variance in estimates of abundance  
 84 than design-based estimates. Model-based inference also enables the use of

85 data from opportunistic surveys, for example, incidental data arising from  
86 “ecotourism” cruises (Williams *et al.*, 2006).

87 Our aims in creating a spatial model of a biological population are usu-  
88 ally two-fold: (i) estimating overall abundance and (ii) investigating the re-  
89 lationship between abundance and environmental covariates. As with any  
90 predictions that are outside the range of the data, one should heed the usual  
91 warnings regarding extrapolation. For example, if a model contains eleva-  
92 tion as a covariate, predictions at high, unsampled elevations are unlikely to  
93 be reliable. Frequently, maps of abundance or density are required and any  
94 spurious predictions can be visually assessed, as well as by plotting a histo-  
95 gram of the predicted values. A sensible definition of the region of interest  
96 avoids prediction outside the range of the data.

97 In this article we review the current state of spatial modelling of detection-  
98 corrected count data, illustrating some recent developments useful to applied  
99 ecologists. The methods discussed have been available in Distance software  
100 (Thomas *et al.*, 2010) for some time but the recent advances covered here  
101 have been implemented in a new R package, `dsm` (Miller *et al.*, 2013) and are  
102 to be incorporated into Distance.

103 Throughout this article a motivating data set is used to illustrate the  
104 methods. These data are sightings of pantropical spotted dolphins (*Stenella*  
105 *attenuata*) during April and May of 1996 in the Gulf of Mexico. Observers  
106 aboard the NOAA vessel Oregon II recorded sightings and environmental co-  
107 variates (see <http://seamap.env.duke.edu/dataset/25> for survey details).  
108 A complete example analysis is provided in Appendix A. The data used in  
109 the analysis are available in the `dsm` package and Distance.

110 The rest of the article reviews approaches for the spatial modelling of  
111 distance sampling data before focussing on the density surface modelling ap-  
112 proach of Hedley & Buckland (2004) to estimate abundance and uncertainty.  
113 We then describe recent advances and provide practical advice regarding  
114 model fitting, formulation and checking. Finally we discuss future directions  
115 for research in spatially modelling detection-corrected count data.

## 116 **Approaches to spatial modelling of distance sampling** 117 **data**

118 Modelling of spatially referenced distance sampling data is equivalent to  
119 modelling spatially-referenced count data, with the additional information  
120 provided by collecting distances to account for imperfect detection. We re-  
121 view recent efforts to model such data; some consist of two steps (correction  
122 for imperfect detection, then spatial modelling), while others jointly estimate  
123 the relevant parameters.

### 124 **TWO-STAGE APPROACHES**

125 The focus of this article is the “count model” of Hedley & Buckland (2004),  
126 we will henceforth refer to this approach as *density surface modelling* (DSM).  
127 Modelling proceeds in two steps: a detection function is fitted to the distance  
128 data to obtain detection probabilities for clusters (flocks, pods, etc.) or in-  
129 dividuals. Counts are allocated to corresponding segments (contiguous tran-  
130 sect sections). A generalised additive model (GAM; e.g. Wood, 2006) is then

constructed with the per-segment counts as the response with either counts or segment areas corrected for detectability (see *Density surface modelling*, below). GAMs provide a flexible class of models that include generalized linear models (GLMs; McCullagh & Nelder, 1989) but extend them with the possible addition of splines to create smooth functions of covariates, random effects terms or correlation structures. We cover advances using this approach in *Recent developments*.

Niemi & Fernández (2010) proposed a Bayesian point process approach to spatial abundance modelling. The density of the objects is described by an intensity function, which included spatially-referenced covariates. Model fitting proceeded in two stages: first the detection function was fitted, then the spatial model (via MCMC) assuming the detection function parameters were known, so detection function uncertainty was not incorporated in the spatial model. A marked point process (Cox & Isham, 1980, Section 5.5) could be used to incorporate cluster size information.

Ver Hoef *et al.* (2013) modeled seal populations in the Bering Sea using a Bayesian spatial model, using a detection function to account for uncertain detection and incorporating additional information from a (frequentist) model of seal haul-outs on ice. The detection function and haul-out model corrected the observed density estimates which were modelled using a Bayesian hierarchical model for the spatial component. The Bayesian hierarchical model was itself split into two parts (*i*) a presence/absence part to allow modelling of the large number of zeros in the data and (*ii*) a density part also used to account for spatial autocorrelation. The analysis shows that when extra information is available (such as telemetry data for the haul-out

156 process) additional insight can be derived.

157 We note that there are many approaches to modelling spatially referenced  
158 count data (Oppel *et al.*, 2011, provides an overview of such methods for  
159 marine bird modelling). Also worthy of note is the approach of Barry &  
160 Welsh (2002) who used a two-stage approach to model presence/absence  
161 then spatial pattern (via two GAMs) to account for zero-inflation.

## 162 ONE-STAGE APPROACHES

163 Rather than fitting two separate models, some authors have combined the  
164 detection function and spatial model fitted (mostly via hierarchical Bayesian  
165 methods). The first of these was Royle *et al.* (2004), who estimated the para-  
166 meters of a specified detection function, formulating an unconditional like-  
167 lihood per-point/line as a function of the unobserved transect abundances.  
168 These unobserved abundances were treated as random effects, integrated out  
169 to give a per-transect likelihood as a function of detection function and ran-  
170 dom effects parameters (linear functions of the environmental covariates).  
171 Due to the multinomial nature of the per-transect likelihood proposed, de-  
172 tection distances must be allocated to bins (e.g. 0-5m, 5-15m, etc). Chelgren  
173 *et al.* (2011) proposed replacing the multinomial per-transect likelihood with  
174 a binomial distribution multiplied by a detection function. The binomial  
175 term collapses the multinomial bins into a single bin and gives the number  
176 of animals detected in the transect, thus allowing the use of exact distances.

177 The work of Schmidt *et al.* (2011) took a similar approach to Royle &  
178 Dorazio (2008), building a presence/absence-type model for clusters, aug-  
179 menting the data with unobserved clusters. The authors then used a Poisson



180 distribution to model cluster size (using a random effect to incorporate over-  
181 dispersion), combining these parts gave a model of individual abundance.  
182 Conn *et al.* (2012) also used a hierarchical Bayesian model but in terms of  
183 abundance rather than density using a super-population/data augmentation  
184 approach (as in Link & Barker, 2009). In their formulation, the whole popu-  
185 lation within the study region was modelled, not just those animals observed  
186 during the survey.

187 Moore & Barlow (2011) adopted a hierarchical Bayesian state-space model,  
188 separating the problem into observation and process components. The pro-  
189 cess component described the underlying population density as it changed  
190 over time and space (though the authors only included strata as a spatial  
191 component). The observation part of the model then linked the process  
192 model to the data via the detection function.

193 Johnson *et al.* (2010) proposed a point process-based model for distance  
194 sampling data. They first assumed that the locations of all individuals in  
195 the survey area (not just those observed) form a realisation of a Poisson pro-  
196 cess. Parameters of the intensity function were then estimated via standard  
197 maximum likelihood methods for point processes (Baddeley & Turner, 2000).  
198 All parameters were estimated jointly so uncertainty from both the spatial  
199 pattern and the detection function was incorporated into variance estimates  
200 of the abundance. This also ensured that correlations between the detection  
201 function and underlying point process were estimated correctly (and did not  
202 falsely inflate or deflate variance estimates). A post-hoc correction factor was  
203 used to address overdispersion unmodelled by spatial covariates (i.e. counts  
204 that do not follow a Poisson mean-variance relationship).

206 Generally very little information is lost by taking a two-stage approach. This  
207 is because transects are typically very narrow compared with the width of the  
208 study area so, provided no significant density variation takes place “across”  
209 the width of the lines or within the point, there is no information in the  
210 distances about the spatial distribution of animals (this is an assumption of  
211 two-stage approaches).

212 Two-stage approaches are effectively “divide and conquer” techniques:  
213 concentrating on the detection function first, and then, given the detection  
214 function, fitting the spatial model. One-stage models are more difficult to  
215 both estimate and check as both steps occur at once; models are potentially  
216 simpler from the perspective of the user and perhaps more mathematically  
217 elegant.

218 Two-stage models have the disadvantage that to accurately quantify model  
219 uncertainty one must appropriately combine uncertainty from the detection  
220 function and spatial models. This can be challenging; however, the alternat-  
221 ive of ignoring uncertainty from the detection process (e.g. Niemi & Fernán-  
222 dez, 2010) can produce confidence or credible intervals for abundance estim-  
223 ates that have coverage below the nominal level. More information regarding  
224 how variance estimation is addressed for DSMs is given in *Recent develop-*  
225 *ments*.

## 226 Density surface modelling

227 This section focuses on modelling the density/abundance estimation stage of  
228 the DSM approach introduced previously. Both line and point transects can  
229 be used, but if lines are used then they are split into contiguous *segments*  
230 (indexed by  $j$ ), which are of length  $l_j$ . Segments should be small enough such  
231 that neither density of objects nor covariate values vary appreciably within  
232 a segment (making the segments approximately square is usually sufficient;  
233  $2w \times 2w$ , where  $w$  is the truncation distance). The area of each segment enters  
234 the model as (or as part of) an offset: the area of segment  $j$  is  $A_j = 2wl_j$   
235 and for point  $j$  is  $A_j = \pi w^2$ .

236 Count or estimated abundance (per segment or point) is then modelled  
237 as a sum of smooth functions of covariates ( $z_{jk}$  with  $k$  indexing the covari-  
238 ates, e.g., location, sea surface temperature, weather conditions; measured at  
239 the segment/point level) using a generalized additive model. Smooth func-  
240 tions are modelled as splines, providing flexible unidimensional (and higher-  
241 dimensional) curves (and surfaces, etc) that describe the relationship between  
242 the covariates and response. Wood (2006) and Ruppert *et al.* (2003) provide  
243 more in-depth introductions to smoothing and generalized additive models.

244 We begin by describing a formulation where only covariates measured  
245 per-segment (e.g. habitat, Beaufort sea state) are included in the detection  
246 function. We later expand this simple formulation to include observation  
247 level covariates (e.g., cluster size, species)

249 The model for the count per segment is:

$$\mathbb{E}(n_j) = \hat{p}_j A_j \exp \left[ \beta_0 + \sum_k f_k(z_{jk}) \right],$$

250 where the  $f_k$ s are smooth functions of the covariates and  $\beta_0$  is an intercept  
 251 term. Multiplying the segment area ( $A_j$ ) by the probability of detection ( $\hat{p}_j$ )  
 252 gives the *effective area* for segment  $j$ . If there are no covariates other than  
 253 distance in the detection function then the probability of detection is constant  
 254 for all segments (i.e.,  $\hat{p}_j = \hat{p}$ ,  $\forall j$ ). The distribution of  $n_j$  can be modelled  
 255 as an overdispersed Poisson, negative binomial, or Tweedie distribution (see  
 256 *Recent developments*).

257 Fig. 2 shows the raw observations of the dolphin data, along with the  
 258 transect lines, overlaid on the depth data. A half-normal detection function  
 259 was fitted to the distances and is shown in Fig. 1. Fig. 3 shows a DSM fitted  
 260 to the dolphin data. The top panel shows predictions from a model where  
 261 depth was the only covariate, the bottom panel shows predictions where  
 262 a (bivariate) smooth of spatial location was also included. Comparing the  
 263 models using GCV score, the latter had a considerably lower score (39.12 vs  
 264 48.46) and so would be selected as our preferred model.

265 As well as simply calculating abundance estimates, relationships between  
 266 covariates and abundance can be illustrated via plots of marginal smooths.  
 267 The effect of depth on abundance (on the scale of the link function) for the  
 268 dolphin data can be seen in Fig. 4.

269 An alternative to modelling counts is to use the per-segment/circle abund-

270 ance using distance sampling estimates as the response. In this case we  
 271 replace  $n_j$  by:

$$\hat{N}_j = \sum_{r=1}^{R_j} \frac{s_{jr}}{\hat{p}_j},$$

272 where  $R_j$  is the number observations in segment  $j$  and  $s_{jr}$  is the size of the  
 273  $r^{\text{th}}$  cluster in segment  $j$  (if the animals occur individually then  $s_{jr} = 1, \forall j, r$ ).

274 The following model is then fitted:

$$\mathbb{E}(\hat{N}_j) = A_j \exp \left[ \beta_0 + \sum_k f_k(z_{jk}) \right],$$

275 where  $\hat{N}_j$ , as with  $n_j$ , is assumed to follow an overdispersed Poisson, negative  
 276 binomial, or Tweedie distribution (see *Recent developments*, below). Note  
 277 that the offset ( $A_j$ ) is now the area of segment/point rather than effective  
 278 area of the segment/point. Although  $\hat{N}_j$  can always be modelled instead of  
 279  $n_j$ , it seems preferable to use  $n_j$  when possible, as one is then modelling actual  
 280 (integer) counts as the response rather than estimates. Note that although  
 281  $\hat{N}_j$  may take non-integer values, this does not present an estimation problem  
 282 for the response distributions covered here.

### 283 *DSM with covariates at the observation level*

284 The above models consider the case where the covariates are measured at  
 285 the segment/point level. Often covariates ( $z_{ij}$ , for individual/cluster  $i$  and  
 286 segment/point  $j$ ) are collected on the level of observations; for example sex  
 287 or cluster size of the observed object or identity of the observer. In this  
 288 case the probability of detection is a function of the object (individual or

289 cluster) level covariates  $\hat{p}(z_i)$ . Object level covariates can be incorporated  
 290 into the model by adopting the following estimator of the per-segment/point  
 291 abundance:

$$\hat{N}_j = \sum_{r=1}^{R_j} \frac{s_{jr}}{\hat{p}(z_{rj})}.$$

292 Density, rather than abundance, can be modelled by excluding the offset  
 293 and instead dividing the count (or estimated abundance) by the area of the  
 294 segment/point (and weighting observations by the segment/point areas). We  
 295 concentrate on abundance here; see Hedley & Buckland (2004) for further  
 296 details on modelling density.

## 297 PREDICTION

298 A DSM can be used to predict abundance over a larger/different area than  
 299 was originally surveyed. In that case the investigator must create a series  
 300 of prediction cells over the prediction region. For each cell the covariates  
 301 included in the DSM must be available; the area of each cell is also required.  
 302 Having made predictions for each cell, these can be plotted as an abundance  
 303 map (as in Fig. 3) and, by summing over cells, an overall estimate of abund-  
 304 ance can be calculated. It is worth noting that using prediction grid cells  
 305 that are smaller than the resolution of the spatially referenced data has no  
 306 effect on abundance/density estimates.

308 Estimating the variance of abundances calculated using a DSM is not straight-  
 309 forward: uncertainty from the estimated parameters of the detection function  
 310 must be incorporated into the spatial model. A second consideration is that  
 311 in a line transect survey, abundances in adjacent segments are likely to be  
 312 correlated; failure to account for this spatial autocorrelation will lead to ar-  
 313 tificially low variance estimates and hence misleadingly narrow confidence  
 314 intervals.

315 Hedley & Buckland (2004) describe a method of calculating the variance  
 316 in the abundance estimates using a parametric bootstrap, resampling from  
 317 the residuals of the fitted model. The bootstrap procedure is as follows.

318 Denote the fitted values for the model to be  $\hat{\boldsymbol{\eta}}$ . For  $b = 1, \dots, B$  (where  
 319  $B$  is the number of resamples required).

- 320 1. Resample (with replacement) the per-segment/point residuals, store  
 321 the values in  $\mathbf{r}_b$ .
- 322 2. Refit the model but with the response set to  $\hat{\boldsymbol{\eta}} + \mathbf{r}_b$  (where  $\hat{\boldsymbol{\eta}}$  are the  
 323 fitted values from the original model).
- 324 3. Take the predicted values for the new model and store them.

325 From the predicted values stored in the last step the variance originating in  
 326 the spatial part of the model can be calculated. The total variance of the  
 327 abundance estimate (over the whole region of interest or sub-areas) can then  
 328 be found by combining the variance estimate from the bootstrap procedure  
 329 with the variance of the probability of detection from the detection function

330 model using the delta method (which assumes that the two components of  
331 the variance are independent; Ver Hoef, 2012).

332     The above procedure assumes that there is no correlation in space between  
333 segments, which are usually contiguous along transects. If many animals are  
334 observed in a particular segment then we might expect there to be high num-  
335 bers in the adjacent segments. A moving block bootstrap (MBB; Efron &  
336 Tibshirani, 1993, Section 8.6) can account for some of this spatial autocor-  
337 relation in the variance estimation. The segments are grouped together into  
338 overlapping blocks (so if the block size is 5, block one is segments 1, ..., 5,  
339 block two is segments 2, ..., 6, and so on). Then, at step (2) above, res-  
340 amples are taken at the block level (rather than individual segments within  
341 a transect). Using MMB will account for correlation between the segments at  
342 scales smaller than the block size, inflating the variances accordingly. Block  
343 size can be selected by plotting an autocorrelogram of the residuals from the  
344 DSM.

345     Both bootstrap procedures can also be modified to take detection function  
346 uncertainty into account. Distances are simulated from the fitted detection  
347 function and then the offset is re-calculated by fitting a detection function  
348 to the simulated distances.

349     Uncertainty can be estimated for a given prediction region by calculat-  
350 ing the appropriate quantiles of the resulting abundance estimates (outlier  
351 removal may be required before quantile calculation). DSM uncertainty can  
352 be visualised via a plot of per-cell coefficient of variation obtained by dividing  
353 the standard error for each cell by its predicted abundance (as in Fig. 5).



## Recent developments

### *GAM uncertainty and variance propagation*

Rather than using a bootstrap, one can use GAM theory to construct uncertainty estimates for DSM abundance estimates. This requires that we use the distribution of the parameters in the GAM to simulate model coefficients, using them to generate replicate abundance estimates (further information can found in Wood, 2006, page 245). Such an approach removes the need to refit the model many times, making variance estimation much faster.

Williams *et al.* (2011) go a step further and incorporate the uncertainty in the estimation of the detection function into the variance of the spatial model, albeit only when segment level covariates are in the DSM. Their procedure is to fit the density surface model with an additional random effect term that characterises the uncertainty in the estimation of the detection function (via the derivatives of the probability of detection,  $\hat{p}$ , with respect to their parameters). Variance estimates of the abundance calculated using standard GAM theory will include uncertainty from the estimation of the detection function. A more complete mathematical explanation of this result is given in Appendix B.

We consider that propagating the uncertainty in this manner to be preferable to the MBB because it is more computationally efficient meaning investigators can easily and quickly estimate variances of complex models. The confidence intervals produced via variance propagation appear comparable (if not narrower) than their bootstrap equivalents, while maintaining good coverage (results of a small simulation study are given in Appendix C).

378 Fig. 5 shows a map of the coefficient of variation for the model which  
379 includes both location and depth covariates. Variance has been calculated  
380 using the variance propagation method.

## 381 EDGE EFFECTS

382 Previous work (Ramsay, 2002; Wang & Ranalli, 2007; Wood *et al.*, 2008;  
383 Scott-Hayward *et al.*, 2013; Miller & Wood, submitted) has highlighted the  
384 need to take care when smoothing over areas with complicated boundaries,  
385 e.g., those with rivers, peninsulae or islands. If two parts of the study area  
386 (either side of a river or inlet, say) are inappropriately linked by the model  
387 (i.e. if the distance between the points is measured as a straight line, rather  
388 than taking into account obstacles) then the boundary feature (river, etc)  
389 can be “smoothed across” so positive abundances are predicted in areas where  
390 animals could not possibly occur. Ensuring that a realistic spatial model has  
391 been fitted to the data is essential for valid inference. The soap film smoother  
392 of Wood *et al.* (2008) is an appealing solution: a bivariate smooth function  
393 of location that can be included in any GAM but that allows for boundary  
394 conditions to be estimated and obeyed for a complex study area. Such an  
395 approach can be helpful when uncertainty is estimated via a bootstrap as  
396 edge effects can also cause large, unrealistic predictions which can plague  
397 other smoothers (Bravington & Hedley, 2009).

398 Even if the study area does not have a complicated boundary, edge effects  
399 can still be problematic. Miller (2012) notes that some smoothers have plane  
400 components that tend to cause the fitted surface to increase unrealistically as  
401 predictions are made further away from the locations of survey effort. This

402 problem can be alleviated by the using a different type of smoother (e.g. a  
403 generalisation of thin plate regression splines called *Duchon splines*).

#### 404     TWEEDIE DISTRIBUTION

405 The Tweedie distribution offers a flexible alternative to the quasi-Poisson and  
406 negative binomial distributions as a response distribution when modelling  
407 count data (Candy, 2004). In particular it is useful when there are a high  
408 proportion of zeros in the data (Shono, 2008; Peel *et al.*, 2012) and avoids  
409 multiple-stage modelling of zero-inflated data (as in Barry & Welsh, 2002).

410 The distribution has three parameters parameters: a mean, dispersion  
411 and a third power parameter, which leads to additional flexibility. The dis-  
412 tribution does not change appreciably when the power parameter is changed  
413 by less than 0.1 and therefore a simple line search over the possible values  
414 for the power parameter is usually a reasonable approach to estimating the  
415 parameter. Mark Bravington (pers. comm.) suggested plotting the square  
416 root of the absolute value of the residuals against fitted values; a “flatter”  
417 plot (points forming a horizontal line) give an indication of a “good” value.  
418 We additionally suggest using the metrics described in the next section for  
419 model selection.

420 Appendix D gives further details about the Tweedie distribution (includ-  
421 ing its probability density function and further references).

## 422 Practical advice

423 A flow diagram of the modelling process for creating a DSM is shown in Fig.  
424 6. The diagram shows which methods are compatible with each other and  
425 what the options are for modelling a particular data set.

426 In our experience, it is sensible to obtain a detection function that fits  
427 the data as well as possible and only begin spatial modelling after a satisfact-  
428 ory detection function has been obtained. Model selection for the detection  
429 function can be performed using AIC and model checking using goodness-of-  
430 fit tests given in Burnham *et al.* (2004, Section 11.11). If animals occur in  
431 clusters rather than individually, bias can be incurred due to the higher visib-  
432 ility of larger clusters. It may then be necessary to include size as a covariate  
433 in the detection function (see Buckland *et al.*, 2001, Section 4.8.2.4). For  
434 some species cluster size may change according to location, Ferguson *et al.*  
435 (2006) use two GAMs (one to model observed clusters and one to model the  
436 cluster size) to deal with spatially-varying cluster size amongst delphinids,  
437 though the authors do not present the variance of the resulting predictions.

438 Smooth terms can be selected using (approximate)  $p$ -values (Wood, 2006,  
439 Section 4.8.5). An additional useful technique for covariate selection is to  
440 use an extra penalty for each term in the GAM allowing smooth terms to  
441 be removed from the model during fitting (illustrated in Appendix A; Wood,  
442 2011). Smoothness selection is performed by generalized cross validation  
443 (GCV) score, unbiased risk estimator (UBRE) or restricted maximum likeli-  
444 hood (REML) score. When model covariates are effectively functions of one  
445 another (e.g. depth could be written as a function of location) GCV and

446 UBRE can suffer from optimisation problems (Wood, 2006, Section 4.5.3)  
447 which can lead to unstable models (Wood, 2011). REML provides a fitting  
448 criteria with a more pronounced optima which avoids some problems with  
449 parameter estimation, though caution should always be taken when deal-  
450 ing with highly correlated covariates. A significant drawback of REML is  
451 that scores cannot be used to compare models with different linear terms or  
452 offsets (Wood, 2011), though the  $p$ -value and additional penalty techniques  
453 described above can be used to select model terms. We highly recommend  
454 the use of standard GAM diagnostic plots; Wood (2006) provides further  
455 practical information on GAM model selection and fitting.

456 In the analysis of the dolphin data we included a smooth of location that  
457 nearly doubles the percentage deviance explained (27.3% to 52.7%). One can  
458 see this when comparing the two plots in Fig. 3 and the plot of the depth  
459 (Fig. 2), the plot of the model containing only a smooth of depth looks very  
460 similar to the raw plot of the depth data. Using a smooth of location can be  
461 a primitive way to account for spatial autocorrelation and/or as a proxy for  
462 other spatially varying covariates that are unavailable.

463 A more sophisticated way to account for spatial autocorrelation between  
464 segments (within transects) is to use an autocorrelation structure within the  
465 DSM (e.g. autoregressive models). Appendix A shows an example using  
466 generalized additive mixed model (GAMMs; Wood, 2006, Section 6.6, see  
467 Appendix A for an example) to construct an autoregressive (lag 1) correla-  
468 tion structure. This gives a significant reduction in variance, tightening the  
469 confidence interval around the abundance estimate.

470 In the analysis presented here, spatial location has been transformed from

471 latitude and longitude to kilometres north and east of the centre of the sur-  
472 vey region at  $(27.01^\circ, -88.3^\circ)$ . This is because the bivariate smoother used  
473 (the thin plate spline; Wood, 2003) is isotropic: there is only one parameter  
474 controlling the smoothness in both directions. Moving one degree in latitude  
475 is not the same as moving one degree in longitude and so using kilometres  
476 from the centre of the study region makes the covariates isotropic. Using  
477 metric units rather than non-standard units of measure such as degrees or  
478 feet throughout makes analysis much easier.

479 A smooth of an environment-level covariate such as depth can be very  
480 useful for assessing the relationships between abundance and the covariate  
481 (as in Fig. 4). Caution should be employed when interpreting smooth re-  
482 lationships and abundance estimates, especially if there are gaps over the  
483 range of covariate values. Large counts may occur at large values of depth  
484 but if no further observations occur at such a large value, then investigators  
485 should be skeptical of any relationship.

## 486 Discussion

487 The use of model-based inference for determining abundance and spatial  
488 distribution from distance sampling data presents new opportunities in the  
489 field of population assessment. Spatial models can be particularly useful  
490 when it comes to prediction: making predictions for some subset of the study  
491 area relies on stratification in design-based methods and as such can be rather  
492 limited. Our models also allow inference from a sample of sightings to a  
493 population in a study area without depending upon a random sample design,

494 and therefore data collected from "platforms of opportunity" (Williams *et al.*,  
495 2006) can be used (although a well designed survey is always preferable).

496 Unbiased estimates are dependent upon either (i) distribution of sampling  
497 effort being random throughout the study area (for design-based inference)  
498 or (ii) model correctness (for model-based inference). It is easier to have  
499 confidence in the former rather than in the latter because our models are  
500 always wrong. Nevertheless model-based inference will play an increasing  
501 role in population assessment as the availability of spatially-referenced data  
502 increases.

503 The field is quickly evolving to allow modelling of more complex data  
504 building on the basic ideas of density surface modelling. We expect to see  
505 large advances in temporal inferences and the handling of zero-inflated data  
506 and spatial correlation. These should become more mainstream as modern  
507 spatio-temporal modelling techniques are adopted. Petersen *et al.* (2011)  
508 provided a very basic framework for temporal modelling; their model included  
509 "before" and "after" smooth terms to quantify the impact of the construction  
510 of an offshore windfarm. Zero-inflation in count data may be problematic  
511 and two-stage approaches such as Barry & Welsh (2002) as well as more flex-  
512 ible response distributions made possible by Rigby & Stasinopoulos (2005)  
513 have yet to be exploited by those using distance sampling data. Spatial  
514 autocorrelation can be accounted for via approaches that explicitly intro-  
515 duce correlations such as generalized estimating equations (GEEs; Hardin &  
516 Hilbe, 2003) or generalized additive mixed models or via mechanisms such  
517 as that of Skaug (2006), which allow observations to cluster according to one  
518 of several states (such as high vs low density patches, possibly in response to

519 temporary agglomerations of prey, although the mechanism is unimportant).  
520 These advances should assist both modellers and wildlife managers to make  
521 optimal conservation decisions.

522 Advances in Bayesian computation (INLA; Rue *et al.*, 2009), make one-  
523 step, Bayesian, density surface models computationally feasible (as INLA  
524 is an alternative to MCMC). An important step toward such models will  
525 be incorporation of detection function estimation into the spatial model.  
526 We anticipate that such a direct modelling technique will dominate future  
527 developments in the field.

528 Density surface modelling allows wildlife managers to make best use of the  
529 available spatial data to understand patterns of abundance, and hence make  
530 better conservation decisions (e.g., about reserve or development placement).  
531 The recent advances mentioned here increase the reliability of the outputs  
532 from a modelling exercise, and hence the efficacy of these decisions. Density  
533 surface modelling from survey data is an active area of research, and we look  
534 forward to further improvements and extensions in the near future.

## 535 Acknowledgments

536 We wish to thank Paul Conn, another anonymous reviewer, and the asso-  
537 ciate editor for their helpful comments. DLM wishes to thank Mark Brav-  
538 ington and Sharon Hedley for their detailed discussions and for providing  
539 code for their variance propagation method. Funding for the implementa-  
540 tion of the recent advances into the `dsm` package and Distance software came  
541 from the US Navy, Chief of Naval Operations (Code N45), grant number



542 N00244-10-1-0057.

## 543 References

- 544 Baddeley, A. & Turner, R. (2000) Practical maximum pseudolikelihood for spatial  
545 point patterns. *Australian & New Zealand Journal of Statistics*, **42**, 283–322.
- 546 Barry, S.C. & Welsh, A.H. (2002) Generalized additive modelling and zero inflated  
547 count data. *Ecological Modelling*, **157**, 179–188.
- 548 Bravington, M.V. & Hedley, S.L. (2009) Antarctic minke whale abundance estim-  
549 ates from the second and third circumpolar IDCR/SOWER surveys using the  
550 SPLINTR model. Paper SC/61/IA14, IWC Scientific Committee.
- 551 Buckland, S.T., Anderson, D.R., Burnham, K.P., Laake, J.L., Borchers, D.L. &  
552 Thomas, L. (2001) *Introduction to Distance Sampling*. Oxford University Press.
- 553 Buckland, S.T., Anderson, D.R., Burnham, K.P., Laake, J.L., Borchers, D.L. &  
554 Thomas, L. (2004) *Advanced Distance Sampling*. Oxford University Press.
- 555 Burnham, K.P., Buckland, S.T., Laake, J.L., Borchers, D.L., Marques, T.A.,  
556 Bishop, J.R. & Thomas, L. (2004) Further topics in distance sampling. *Ad-  
557 vanced Distance Sampling* (eds. S.T. Buckland, D.R. anderson, K.P. Burnham,  
558 J.L. Laake, D.L. Borchers & L. Thomas). Oxford University Press.
- 559 Candy, S. (2004) Modelling catch and effort data using generalised linear models,  
560 the Tweedie distribution, random vessel effects and random stratum-by-year  
561 effects. *CCAMLR Science*, **11**, 59–80.
- 562 Chelgren, N.D., Samora, B., Adams, M.J. & McCreary, B. (2011) Using spati-  
563 otemporal models and distance sampling to map the space use and abundance  
564 of newly metamorphosed western toads (*Anaxyrus boreas*). *Herpetological Con-  
565 servation and Biology*, **6**, 175–190.
- 566 Conn, P.B., Laake, J.L. & Johnson, D.S. (2012) A hierarchical modeling framework  
567 for multiple observer transect surveys. *PLoS ONE*, **7**, e42294.
- 568 Cox, D.R. & Isham, V. (1980) *Point Processes*. Monographs on Applied Probability  
569 and Statistics. Chapman and Hall. ISBN 9780412219108.
- 570 Efron, B. & Tibshirani, R.J. (1993) *An Introduction to the Bootstrap*. Chapman &  
571 Hall/CRC. ISBN 9780412042317.
- 572 Ferguson, M.C., Barlow, J., Fiedler, P., Reilly, S.B. & Gerrodette, T. (2006) Spatial  
573 models of delphinid (family Delphinidae) encounter rate and group size in the  
574 eastern tropical Pacific Ocean. *Ecological Modelling*, **193**, 645–662.

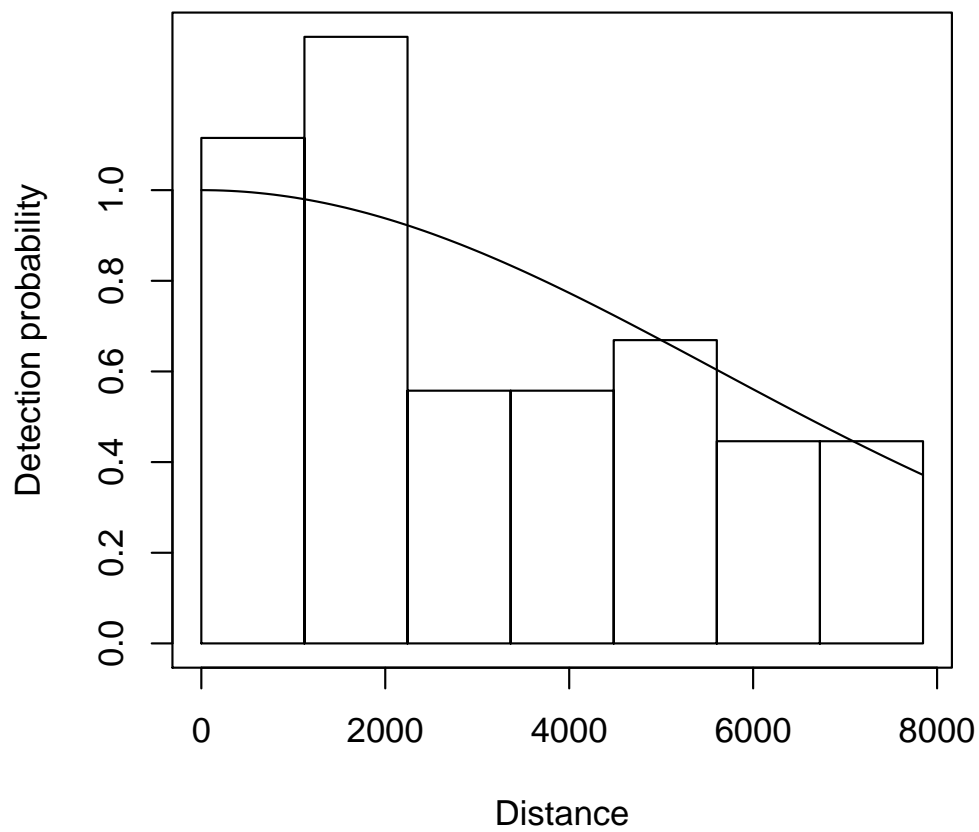
- 575 Halpin, P., Read, A., Fujioka, E., Best, B., Donnelly, B., Hazen, L., Kot, C.,  
576 Urian, K., LaBrecque, E., Dimatteo, A., Cleary, J., Good, C., Crowder, L.  
577 & Hyrenbach, K.D. (2009) OBIS-SEAMAP: The world data center for marine  
578 mammal, sea bird, and sea turtle distributions. *Oceanography*, **22**, 104–115.
- 579 Hardin, J. & Hilbe, J. (2003) Generalized Estimating Equations. Chapman and  
580 Hall/CRC, London, UK.
- 581 Hedley, S.L. & Buckland, S.T. (2004) Spatial models for line transect sampling.  
582 *Journal of Agricultural, Biological, and Environmental Statistics*, **9**, 181–199.
- 583 Johnson, D.S., Laake, J.L. & Ver Hoef, J.M. (2010) A model-based approach for  
584 making ecological inference from distance sampling data. *Biometrics*, **66**, 310–  
585 318.
- 586 Link, W.A. & Barker, R.J. (2009) *Bayesian Inference: with ecological applications*.  
587 Academic Press, London, UK.
- 588 Marques, T.A., Thomas, L., Fancy, S. & Buckland, S.T. (2007) Improving estimates  
589 of bird density using multiple-covariate distance sampling. *The Auk*, **124**, 1229–  
590 1243.
- 591 McCullagh, P. & Nelder, J.A. (1989) *Generalized Linear Models*. Chapman &  
592 Hall/CRC.
- 593 Miller, D.L. (2012) *On smooth models for complex domains and distances*. Ph.D.  
594 thesis, University of Bath.
- 595 Miller, D.L., Rexstad, E.A., Burt, M.L., Bravington, M.V. & Hedley, S.L. (2013)  
596 *dsm: Density surface modelling of distance sampling data*.  
597 URL <http://github.com/dill/dsm>
- 598 Miller, D.L. & Wood, S.N. (submitted) Finite area smoothing with generalized  
599 distance splines.
- 600 Moore, J.E. & Barlow, J. (2011) Bayesian state-space model of fin whale abundance  
601 trends from a 1991-2008 time series of line-transect surveys in the California  
602 Current. *Journal of Applied Ecology*, **48**, 1195–1205.
- 603 Niemi, A. & Fernández, C. (2010) Bayesian spatial point process modeling of line  
604 transect data. *Journal of Agricultural, Biological, and Environmental Statistics*,  
605 **15**, 327–345.
- 606 Oppel, S., Meirinho, A., Ramírez, I., Gardner, B., O’Connell, A., Miller, P. &  
607 Louzao, M. (2011) Comparison of five modelling techniques to predict the spatial  
608 distribution and abundance of seabirds. *Biological Conservation*, **156**, 94–104.

- 609 Peel, D., Bravington, M.V., Kelly, N., Wood, S.N. & Knuckey, I. (2012) A Model-  
610 Based Approach to Designing a Fishery-Independent Survey. *Journal of Agri-  
611 cultural, Biological, and Environmental Statistics*, **18**, 1–21.
- 612 Petersen, I.K., MacKenzie, M.L., Rexstad, E.A., Wisz, M.S. & Fox, A.D. (2011)  
613 Comparing pre- and post-construction distributions of long-tailed ducks *Clan-  
614 gula hyemalis* in and around the Nysted offshore wind farm, Denmark: a quasi-  
615 designed experiment accounting for imperfect detection, local surface features  
616 and autocorrelation. Technical report 2011-1, Centre for Research into Environ-  
617 mental and Ecological Modelling.
- 618 Ramsay, T. (2002) Spline smoothing over difficult regions. *Journal of the Royal  
619 Statistical Society. Series B, Statistical Methodology*, **64**, 307–319.
- 620 Rigby, R. & Stasinopoulos, D. (2005) Generalized additive models for location, scale  
621 and shape. *Journal of the Royal Statistical Society-Series C Applied Statistics*,  
622 **54**, 507–554.
- 623 Royle, J.A. & Dorazio, R.M. (2008) *Hierarchical Modeling and Inference in Ecology*.  
624 Academic Press, London, UK.
- 625 Royle, J., Dawson, D. & Bates, S. (2004) Modeling abundance effects in distance  
626 sampling. *Ecology*, **85**, 1591–1597.
- 627 Rue, H., Martino, S. & Chopin, N. (2009) Approximate Bayesian inference for  
628 latent Gaussian models by using integrated nested Laplace approximations. *J.  
629 R. Statist. Soc. B*, **71**, 319–392.
- 630 Ruppert, D., Wand, M. & Carroll, R.J. (2003) *Semiparametric Regression*. Cam-  
631 bridge Series on Statistical and Probabilistic Mathematics. Cambridge University  
632 Press.
- 633 Schmidt, J.H., Rattenbury, K.L., Lawler, J.P. & Maccluskie, M.C. (2011) Using  
634 distance sampling and hierarchical models to improve estimates of Dall’s sheep  
635 abundance. *The Journal of Wildlife Management*, **76**, 317–327.
- 636 Scott-Hayward, L.A.S., MacKenzie, M.L., Donovan, C.R., Walker, C.G. & Ashe,  
637 E. (2013) Complex region spatial smoother (CReSS). *Journal of Computational  
638 and Graphical Statistics*.
- 639 Shono, H. (2008) Application of the Tweedie distribution to zero-catch data in  
640 CPUE analysis. *Fisheries Research*, **93**, 154–162.
- 641 Skaug, H.J. (2006) Markov modulated Poisson processes for clustered line transect  
642 data. *Environmental and Ecological Statistics*, **13**, 199–211.

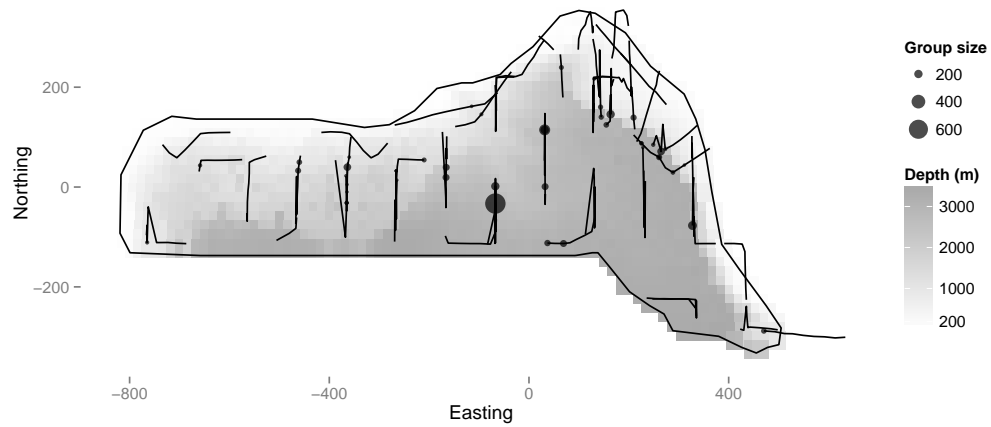
- 643 Thomas, L., Buckland, S.T., Rexstad, E.A., Laake, J.L., Strindberg, S., Hedley,  
644 S.L., Bishop, J.R., Marques, T.A. & Burnham, K.P. (2010) Distance software:  
645 design and analysis of distance sampling surveys for estimating population size.  
646 *Journal of Applied Ecology*, **47**, 5–14.
- 647 Ver Hoef, J.M. (2012) Who invented the delta method? *The American Statistician*,  
648 **66**, 124–127.
- 649 Ver Hoef, J.M., Cameron, M.F., Boveng, P.L., London, J.M. & Moreland, E.E.  
650 (2013) A spatial hierarchical model for abundance of three ice-associated seal  
651 species in the eastern Bering Sea. *Statistical Methodology*, pp. 1–44.
- 652 Wang, H. & Ranalli, M. (2007) Low-rank smoothing splines on complicated do-  
653 mains. *Biometrics*, **63**, 209–217.
- 654 Williams, R., Hedley, S.L., Branch, T.A., Bravington, M.V., Zerbini, A.N. & Find-  
655 lay, K.P. (2011) Chilean blue whales as a case study to illustrate methods to  
656 estimate abundance and evaluate conservation status of rare species. *Conserva-  
657 tion Biology*, **25**, 526–535.
- 658 Williams, R., Hedley, S.L. & Hammond, P. (2006) Modeling distribution and  
659 abundance of Antarctic baleen whales using ships of opportunity. *Ecology and  
660 Society*, **11**, 1.
- 661 Wood, S.N. (2003) Thin plate regression splines. *Journal of the Royal Statistical  
662 Society. Series B, Statistical Methodology*, **65**, 95–114.
- 663 Wood, S.N. (2006) *Generalized Additive Models: An introduction with R*. Chapman  
664 & Hall/CRC.
- 665 Wood, S.N. (2011) Fast stable restricted maximum likelihood and marginal like-  
666 lihood estimation of semiparametric generalized linear models. *Journal of the  
667 Royal Statistical Society. Series B, Statistical Methodology*, **73**, 3–36.
- 668 Wood, S.N., Bravington, M.V. & Hedley, S.L. (2008) Soap film smoothing. *Journal  
669 of the Royal Statistical Society. Series B, Statistical Methodology*, **70**, 931–955.

## Figures

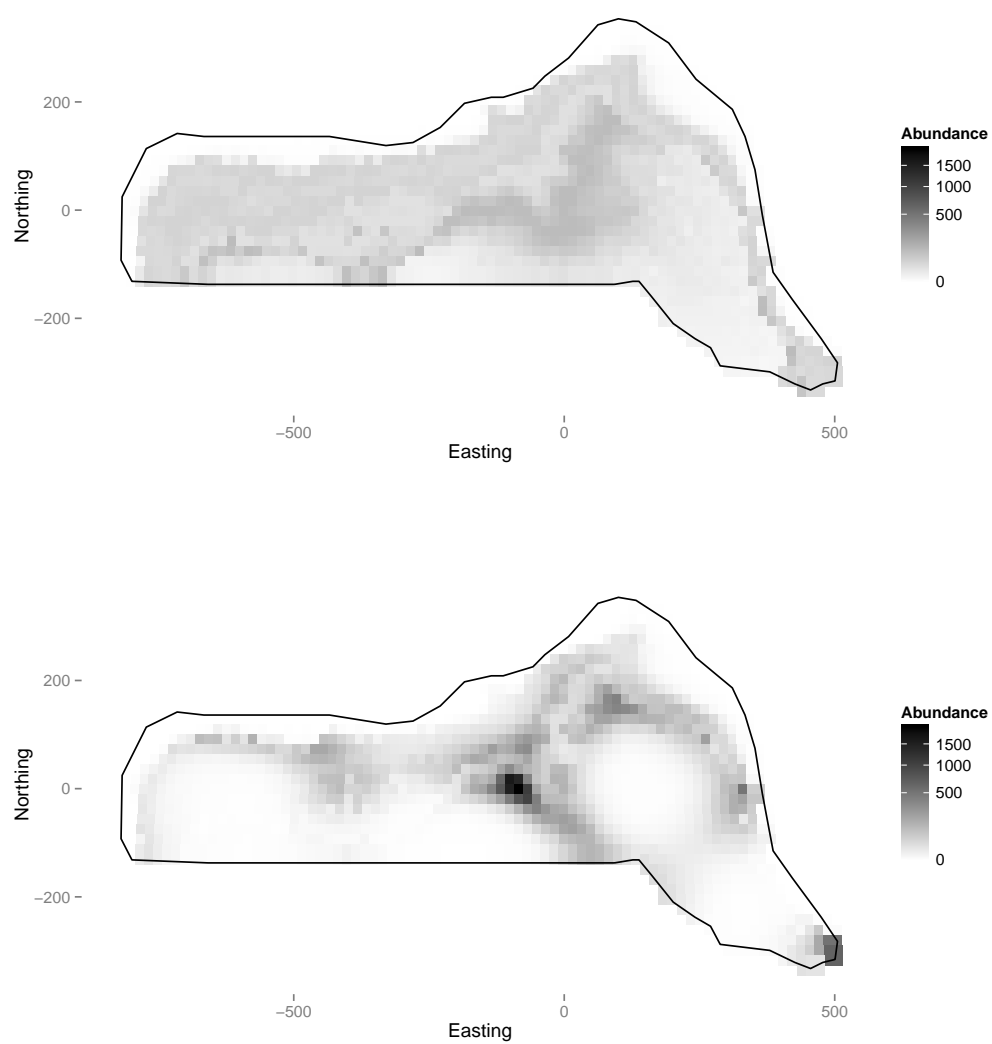
**Fig. 1** Estimated detection function for pantropical dolphin clusters overlaid onto the scaled histogram of observed distances. Distances are recorded in metres.



**Fig. 2** The region, transect centrelines and location of detected pantropical dolphin clusters, where size of circle corresponds to the cluster size, overlaid onto depth data.

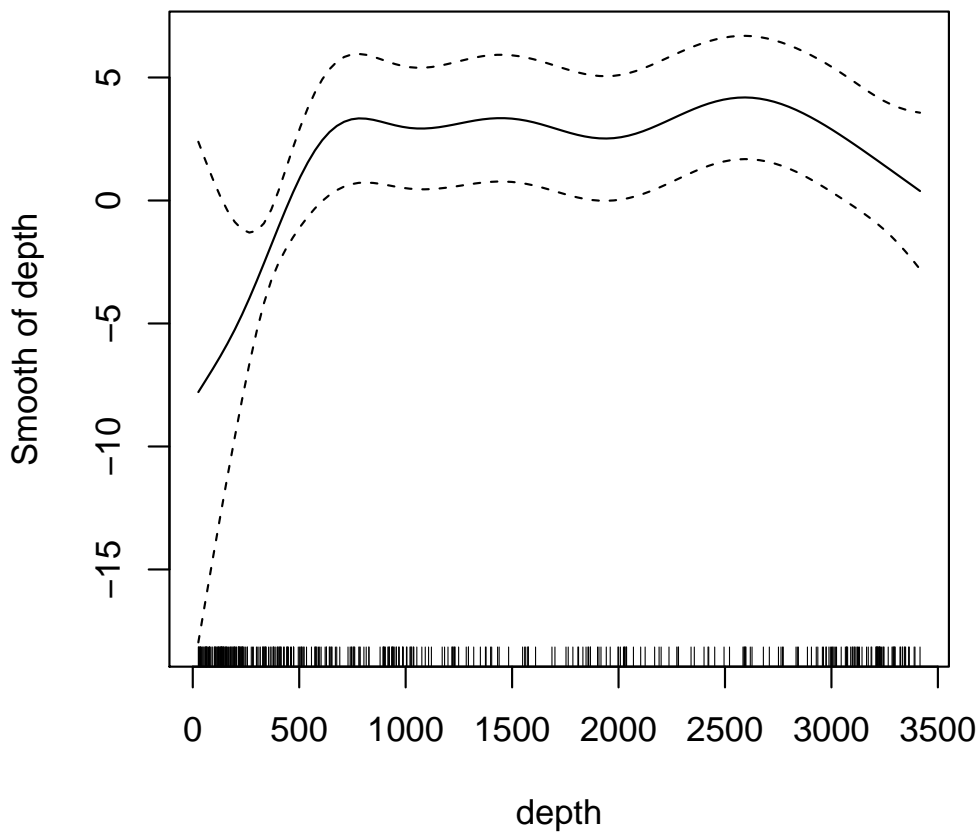


**Fig. 3** Predicted abundance of dolphins from the DSM using only depth as an explanatory variable (top) and the model using both depth and location (bottom).

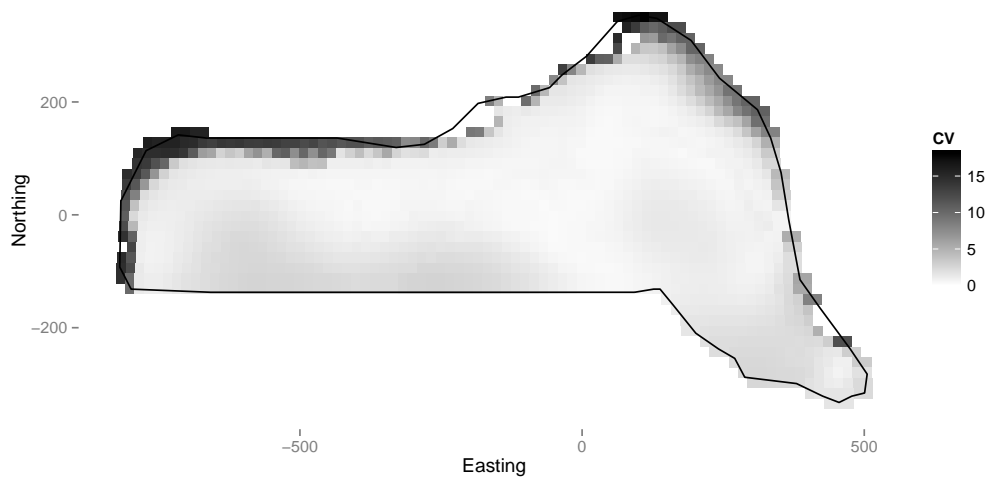




**Fig. 4** Plot of the effect on the response of depth, given location (from the model with both depth and location smooths). Note that it is possible to draw a straight line between 750m and 3000m within the confidence band (between the dashed lines), so the wiggles in the smooth may not be indicative of any relationship. What is clear is that there the estimated number of dolphins increases up to a water depth of about 500m. The rug ticks at the bottom of the plot indicate we have good coverage of the range of depth values in the survey area. Note that the  $y$  axis in such plots is on the scale of the link function (log in this case), so care should be taken in their interpretation.



**Fig. 5** Map of the coefficients of variation for the model with smooths of both depth and location. Uncertainty was estimated using the variance propagation method of Williams *et al.* (2011). As might be expected, there is high uncertainty where there is low sampling effort (Fig. 2).



**Fig. 6** Flow diagram showing the modelling process for creating a density surface model.

