

1 **Running title:** Spatial models for distance sampling
2 **Number of words:** ~5104
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^{1*}, M. Louise Burt²,**
9 **Eric A. Rexstad², Len Thomas².**

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

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 “spatial model” to refer to any model
50 that includes any spatially referenced covariates, not only those models that
51 include explicit location terms. This article is concerned with combining
52 spatial modelling techniques with distance sampling (Buckland *et al.*, 2001,
53 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
129 individuals. Counts are then summarised per segment (contiguous transect
130 section). A generalised additive model (GAM; e.g. Wood, 2006) is then

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

138 As with the DSM approach, Niemi & Fernández (2010) used a two-step
139 procedure: first fitting a detection function, then using a Bayesian point
140 process to model spatial pattern (fitted using MCMC). Object density was
141 described by an intensity function, which included spatially-referenced co-
142 variates. A possible disadvantage of their approach was that the distance
143 function was assumed fixed once its parameters are estimated, and thus un-
144 certainty may not be correctly propagated into final abundance estimates.

145 Ver Hoef *et al.* (2013) also included separate density and detection mod-
146 els for seals in the Bering sea. However, they were able to separate the
147 detection process into three components: (i) incomplete detection on the
148 transect line, (ii) declining detection probability as a function of distance,
149 and (iii) availability bias (as seals could only be observed when hauled out
150 on ice flows). After correcting counts for uncertain detection, they used a
151 hierarchical, zero-inflated spatial regression model to estimate abundance,
152 propagating variance associated with each stage of modelling into final es-
153 timates. The analysis shows that when extra information is available (such
154 as telemetry data for the haul-out process) additional insight can be derived.

155 We note that there are many approaches to modelling spatially referenced

156 count data (Oppel *et al.*, 2011, provides an overview of such methods for
157 marine bird modelling). Also worthy of note is the approach of Barry &
158 Welsh (2002) who used a two-stage approach to model presence/absence then
159 spatial distribution (each via a separate GAM) to account for zero-inflation.

160 ONE-STAGE APPROACHES

161 Rather than fitting two separate models, some authors have estimated para-
162 meters of the detection and spatial models simultaneously. Perhaps the first
163 such example was Royle *et al.* (2004), who considered an integrated likeli-
164 hood model for point and line transects. The approach views abundance as
165 a nuisance variable which was integrated out of the likelihood, but inferences
166 may still be made about factors affecting underlying density (including co-
167 variate effects). This approach was originally developed for binned distance
168 data, but was extended by Chelgren *et al.* (2011) for continuous distance
169 data.

170 Both Schmidt *et al.* (2011) and Conn *et al.* (2012) took data augmentation
171 approaches to add unobserved clusters within their hierarchical Bayesian
172 models. Schmidt *et al.* (2011) used a presence/absence-type model and a
173 super-population approach (as in Royle & Dorazio, 2008). Conn *et al.* (2012)
174 augmented observations only within the sampled transects using RJMCMC.
175 Looking at the problem with at a coarser spatial resolution (stratum-level),
176 Moore & Barlow (2011) separated the problem into observation and process
177 components using a state-space model. The process component described
178 the underlying population density as it changed over time and space, which
179 was linked to the data via the detection function.

180 Another point process-based approach is that of Johnson *et al.* (2010),
181 who used a Poisson process to model the locations of individuals in the survey
182 area. Unlike Niemi & Fernández (2010), parameters of the intensity func-
183 tion were estimated jointly with detection function parameters via standard
184 maximum likelihood methods for point processes (Baddeley & Turner, 2000)
185 (allowing uncertainty from both the spatial pattern and detection function
186 to be included in variance estimates). A post-hoc correction factor was used
187 to address overdispersion unmodelled by spatial covariates (i.e. counts that
188 do not follow a Poisson mean-variance relationship).

189 ONE- VS. TWO-STAGE APPROACHES

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

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

202 Two-stage models have the disadvantage that to accurately quantify model
203 uncertainty one must appropriately combine uncertainty from the detection

204 function and spatial models. This can be challenging; however, the alternat-
205 ive of ignoring uncertainty from the detection process (e.g. Niemi & Fernán-
206 dez, 2010) can produce confidence or credible intervals for abundance estim-
207 ates that have coverage below the nominal level. More information regarding
208 how variance estimation is addressed for DSMs is given in *Recent develop-*
209 *ments*.

210 Density surface modelling

211 This section focuses on modelling the density/abundance estimation stage of
212 the DSM approach introduced previously. Both line and point transects can
213 be used, but if lines are used then they are split into contiguous *segments*
214 (indexed by j), which are of length l_j . Segments should be small enough such
215 that neither density of objects nor covariate values vary appreciably within
216 a segment (making the segments approximately square is usually sufficient;
217 $2w \times 2w$, where w is the truncation distance). The area of each segment enters
218 the model as (or as part of) an offset: the area of segment j is $A_j = 2wl_j$
219 and for point j is $A_j = \pi w^2$.

220 Count or estimated abundance (per segment or point) is then modelled
221 as a sum of smooth functions of covariates (z_{jk} with k indexing the covari-
222 ates, e.g., location, sea surface temperature, weather conditions; measured at
223 the segment/point level) using a generalized additive model. Smooth func-
224 tions are modelled as splines, providing flexible unidimensional (and higher-
225 dimensional) curves (and surfaces, etc) that describe the relationship between
226 the covariates and response. Wood (2006) and Ruppert *et al.* (2003) provide

227 more in-depth introductions to smoothing and generalized additive models.

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

232 COUNT AS RESPONSE

233 The model for the count per segment is:

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

234 where the f_k s are smooth functions of the covariates and β_0 is an intercept
235 term. Multiplying the segment area (A_j) by the probability of detection (\hat{p}_j)
236 gives the *effective area* for segment j . If there are no covariates other than
237 distance in the detection function then the probability of detection is constant
238 for all segments (i.e., $\hat{p}_j = \hat{p}, \forall j$). The distribution of n_j can be modelled
239 as an overdispersed Poisson, negative binomial, or Tweedie distribution (see
240 *Recent developments*).

241 Fig. 2 shows the raw observations of the dolphin data, along with the
242 transect lines, overlaid on the depth data. A half-normal detection function
243 was fitted to the distances and is shown in Fig. 1. Fig. 3 shows a DSM fitted
244 to the dolphin data. The top panel shows predictions from a model where
245 depth was the only covariate, the bottom panel shows predictions where
246 a (bivariate) smooth of spatial location was also included. Comparing the
247 models using GCV score, the latter had a considerably lower score (39.12 vs

248 48.46) and so would be selected as our preferred model.

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

253 An alternative to modelling counts is to use the per-segment/circle abund-
254 ance using distance sampling estimates as the response. In this case we
255 replace n_j by:

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

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

258 The following model is then fitted:

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

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

267 *DSM with covariates at the observation level*

268 The above models consider the case where the covariates are measured at
269 the segment/point level. Often covariates (z_{ij} , for individual/cluster i and
270 segment/point j) are collected on the level of observations; for example sex
271 or cluster size of the observed object or identity of the observer. In this
272 case the probability of detection is a function of the object (individual or
273 cluster) level covariates $\hat{p}(z_i)$. Object level covariates can be incorporated
274 into the model by adopting the following estimator of the per-segment/point
275 abundance:

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

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

281 PREDICTION

282 A DSM can be used to predict abundance over a larger/different area than
283 was originally surveyed. In that case the investigator must create a series
284 of prediction cells over the prediction region. For each cell the covariates
285 included in the DSM must be available; the area of each cell is also required.
286 Having made predictions for each cell, these can be plotted as an abundance
287 map (as in Fig. 3) and, by summing over cells, an overall estimate of abund-

288 ance can be calculated. It is worth noting that using prediction grid cells
289 that are smaller than the resolution of the spatially referenced data has no
290 effect on abundance/density estimates.

291 VARIANCE ESTIMATION

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

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

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

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

309 From the predicted values stored in the last step the variance originating in
310 the spatial part of the model can be calculated. The total variance of the

311 abundance estimate (over the whole region of interest or sub-areas) can then
312 be found by combining the variance estimate from the bootstrap procedure
313 with the variance of the probability of detection from the detection function
314 model using the delta method (which assumes that the two components of
315 the variance are independent; Ver Hoef, 2012).

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

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

333 Uncertainty can be estimated for a given prediction region by calculat-
334 ing the appropriate quantiles of the resulting abundance estimates (outlier
335 removal may be required before quantile calculation). DSM uncertainty can

336 be visualised via a plot of per-cell coefficient of variation obtained by dividing
337 the standard error for each cell by its predicted abundance (as in Fig. 5).

338 Recent developments

339 *GAM uncertainty and variance propagation*

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

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

356 We consider that propagating the uncertainty in this manner to be prefer-
357 able to the MBB because it is more computationally efficient meaning invest-
358 igators can easily and quickly estimate variances of complex models. The

359 confidence intervals produced via variance propagation appear comparable
360 (if not narrower) than their bootstrap equivalents, while maintaining good
361 coverage (results of a small simulation study are given in Appendix C).

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

365 EDGE EFFECTS

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

382 Even if the study area does not have a complicated boundary, edge effects

383 can still be problematic. Miller (2012) notes that some smoothers have plane
384 components that tend to cause the fitted surface to increase unrealistically as
385 predictions are made further away from the locations of survey effort. This
386 problem can be alleviated by the using a different type of smoother (e.g. a
387 generalisation of thin plate regression splines called *Duchon splines*).

388 TWEEDIE DISTRIBUTION

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

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

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

406 Practical advice

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

410 In our experience, it is sensible to obtain a detection function that fits
411 the data as well as possible and only begin spatial modelling after a satisfact-
412 ory detection function has been obtained. Model selection for the detection
413 function can be performed using AIC and model checking using goodness-of-
414 fit tests given in Burnham *et al.* (2004, Section 11.11). If animals occur in
415 clusters rather than individually, bias can be incurred due to the higher visib-
416 ility of larger clusters. It may then be necessary to include size as a covariate
417 in the detection function (see Buckland *et al.*, 2001, Section 4.8.2.4). For
418 some species cluster size may change according to location, Ferguson *et al.*
419 (2006) use two GAMs (one to model observed clusters and one to model the
420 cluster size) to deal with spatially-varying cluster size amongst delphinids,
421 though the authors do not present the variance of the resulting predictions.

422 Smooth terms can be selected using (approximate) p -values (Wood, 2006,
423 Section 4.8.5). An additional useful technique for covariate selection is to
424 use an extra penalty for each term in the GAM allowing smooth terms to
425 be removed from the model during fitting (illustrated in Appendix A; Wood,
426 2011). Smoothness selection is performed by generalized cross validation
427 (GCV) score, unbiased risk estimator (UBRE) or restricted maximum likeli-
428 hood (REML) score. When model covariates are effectively functions of one
429 another (e.g. depth could be written as a function of location) GCV and

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

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

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

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

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

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

470 Discussion

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

478 and therefore data collected from "platforms of opportunity" (Williams *et al.*,
479 2006) can be used (although a well designed survey is always preferable).

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

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

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

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

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

519 Acknowledgments

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

526 N00244-10-1-0057.

527 References

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

- 559 Halpin, P., Read, A., Fujioka, E., Best, B., Donnelly, B., Hazen, L., Kot, C.,
560 Urian, K., LaBrecque, E., Dimatteo, A., Cleary, J., Good, C., Crowder, L.
561 & Hyrenbach, K.D. (2009) OBIS-SEAMAP: The world data center for marine
562 mammal, sea bird, and sea turtle distributions. *Oceanography*, **22**, 104–115.
- 563 Hardin, J. & Hilbe, J. (2003) Generalized Estimating Equations. Chapman and
564 Hall/CRC, London, UK.
- 565 Hedley, S.L. & Buckland, S.T. (2004) Spatial models for line transect sampling.
566 *Journal of Agricultural, Biological, and Environmental Statistics*, **9**, 181–199.
- 567 Johnson, D.S., Laake, J.L. & Ver Hoef, J.M. (2010) A model-based approach for
568 making ecological inference from distance sampling data. *Biometrics*, **66**, 310–
569 318.
- 570 Link, W.A. & Barker, R.J. (2009) *Bayesian Inference: with ecological applications*.
571 Academic Press, London, UK.
- 572 Marques, T.A., Thomas, L., Fancy, S. & Buckland, S.T. (2007) Improving estimates
573 of bird density using multiple-covariate distance sampling. *The Auk*, **124**, 1229–
574 1243.
- 575 McCullagh, P. & Nelder, J.A. (1989) *Generalized Linear Models*. Chapman &
576 Hall/CRC.
- 577 Miller, D.L. (2012) *On smooth models for complex domains and distances*. Ph.D.
578 thesis, University of Bath.
- 579 Miller, D.L., Rexstad, E.A., Burt, M.L., Bravington, M.V. & Hedley, S.L. (2013)
580 *dsm: Density surface modelling of distance sampling data*.
581 URL <http://github.com/dill/dsm>
- 582 Miller, D.L. & Wood, S.N. (submitted) Finite area smoothing with generalized
583 distance splines.
- 584 Moore, J.E. & Barlow, J. (2011) Bayesian state-space model of fin whale abundance
585 trends from a 1991-2008 time series of line-transect surveys in the California
586 Current. *Journal of Applied Ecology*, **48**, 1195–1205.
- 587 Niemi, A. & Fernández, C. (2010) Bayesian spatial point process modeling of line
588 transect data. *Journal of Agricultural, Biological, and Environmental Statistics*,
589 **15**, 327–345.
- 590 Oppel, S., Meirinho, A., Ramírez, I., Gardner, B., O’Connell, A., Miller, P. &
591 Louzao, M. (2011) Comparison of five modelling techniques to predict the spatial
592 distribution and abundance of seabirds. *Biological Conservation*, **156**, 94–104.

- 593 Peel, D., Bravington, M.V., Kelly, N., Wood, S.N. & Knuckey, I. (2012) A Model-
594 Based Approach to Designing a Fishery-Independent Survey. *Journal of Agri-
595 cultural, Biological, and Environmental Statistics*, **18**, 1–21.
- 596 Petersen, I.K., MacKenzie, M.L., Rexstad, E.A., Wisz, M.S. & Fox, A.D. (2011)
597 Comparing pre- and post-construction distributions of long-tailed ducks *Clan-
598 gula hyemalis* in and around the Nysted offshore wind farm, Denmark: a quasi-
599 designed experiment accounting for imperfect detection, local surface features
600 and autocorrelation. Technical report 2011-1, Centre for Research into Environ-
601 mental and Ecological Modelling.
- 602 Ramsay, T. (2002) Spline smoothing over difficult regions. *Journal of the Royal
603 Statistical Society. Series B, Statistical Methodology*, **64**, 307–319.
- 604 Rigby, R. & Stasinopoulos, D. (2005) Generalized additive models for location, scale
605 and shape. *Journal of the Royal Statistical Society-Series C Applied Statistics*,
606 **54**, 507–554.
- 607 Royle, J.A. & Dorazio, R.M. (2008) *Hierarchical Modeling and Inference in Ecology*.
608 Academic Press, London, UK.
- 609 Royle, J., Dawson, D. & Bates, S. (2004) Modeling abundance effects in distance
610 sampling. *Ecology*, **85**, 1591–1597.
- 611 Rue, H., Martino, S. & Chopin, N. (2009) Approximate Bayesian inference for
612 latent Gaussian models by using integrated nested Laplace approximations. *J.
613 R. Statist. Soc. B*, **71**, 319–392.
- 614 Ruppert, D., Wand, M. & Carroll, R.J. (2003) *Semiparametric Regression*. Cam-
615 bridge Series on Statistical and Probabilistic Mathematics. Cambridge University
616 Press.
- 617 Schmidt, J.H., Rattenbury, K.L., Lawler, J.P. & Maccluskie, M.C. (2011) Using
618 distance sampling and hierarchical models to improve estimates of Dall’s sheep
619 abundance. *The Journal of Wildlife Management*, **76**, 317–327.
- 620 Scott-Hayward, L.A.S., MacKenzie, M.L., Donovan, C.R., Walker, C.G. & Ashe,
621 E. (2013) Complex region spatial smoother (CReSS). *Journal of Computational
622 and Graphical Statistics*.
- 623 Shono, H. (2008) Application of the Tweedie distribution to zero-catch data in
624 CPUE analysis. *Fisheries Research*, **93**, 154–162.
- 625 Skaug, H.J. (2006) Markov modulated Poisson processes for clustered line transect
626 data. *Environmental and Ecological Statistics*, **13**, 199–211.

- 627 Thomas, L., Buckland, S.T., Rexstad, E.A., Laake, J.L., Strindberg, S., Hedley,
628 S.L., Bishop, J.R., Marques, T.A. & Burnham, K.P. (2010) Distance software:
629 design and analysis of distance sampling surveys for estimating population size.
630 *Journal of Applied Ecology*, **47**, 5–14.
- 631 Ver Hoef, J.M. (2012) Who invented the delta method? *The American Statistician*,
632 **66**, 124–127.
- 633 Ver Hoef, J.M., Cameron, M.F., Boveng, P.L., London, J.M. & Moreland, E.E.
634 (2013) A spatial hierarchical model for abundance of three ice-associated seal
635 species in the eastern Bering Sea. *Statistical Methodology*, pp. 1–44.
- 636 Wang, H. & Ranalli, M. (2007) Low-rank smoothing splines on complicated do-
637 mains. *Biometrics*, **63**, 209–217.
- 638 Williams, R., Hedley, S.L., Branch, T.A., Bravington, M.V., Zerbini, A.N. & Find-
639 lay, K.P. (2011) Chilean blue whales as a case study to illustrate methods to
640 estimate abundance and evaluate conservation status of rare species. *Conserva-
641 tion Biology*, **25**, 526–535.
- 642 Williams, R., Hedley, S.L. & Hammond, P. (2006) Modeling distribution and
643 abundance of Antarctic baleen whales using ships of opportunity. *Ecology and
644 Society*, **11**, 1.
- 645 Wood, S.N. (2003) Thin plate regression splines. *Journal of the Royal Statistical
646 Society. Series B, Statistical Methodology*, **65**, 95–114.
- 647 Wood, S.N. (2006) *Generalized Additive Models: An introduction with R*. Chapman
648 & Hall/CRC.
- 649 Wood, S.N. (2011) Fast stable restricted maximum likelihood and marginal like-
650 lihood estimation of semiparametric generalized linear models. *Journal of the
651 Royal Statistical Society. Series B, Statistical Methodology*, **73**, 3–36.
- 652 Wood, S.N., Bravington, M.V. & Hedley, S.L. (2008) Soap film smoothing. *Journal
653 of the Royal Statistical Society. Series B, Statistical Methodology*, **70**, 931–955.

654 **Figures**

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

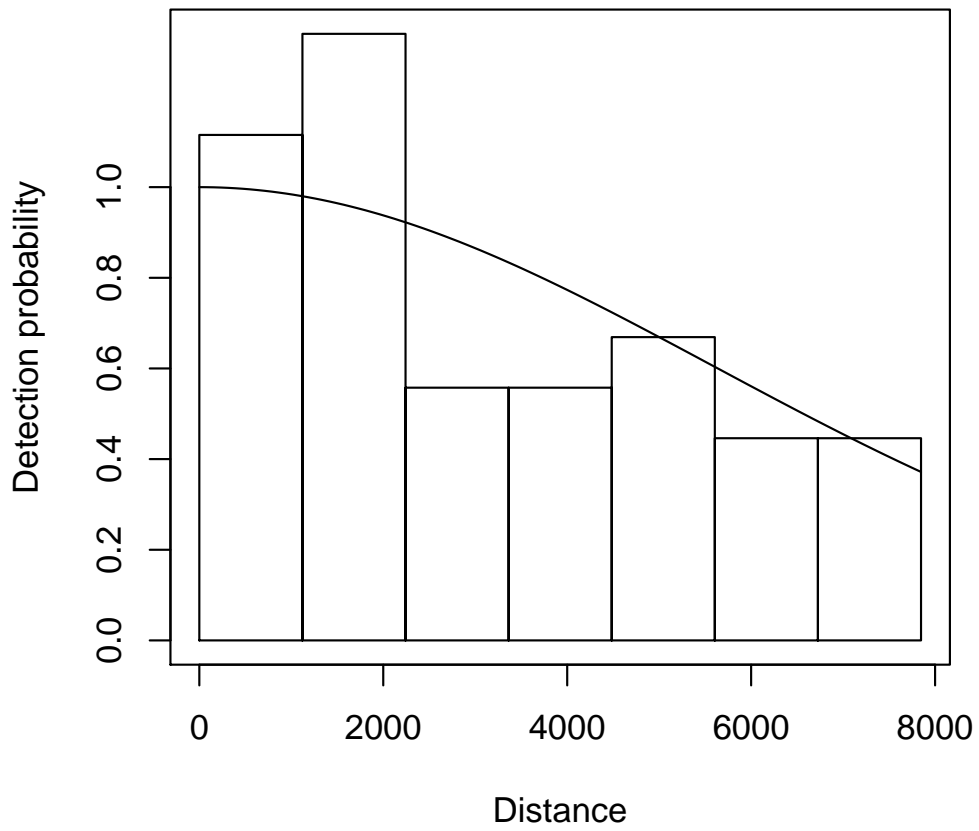


Fig. 2 The region, transect centrelines and location of detected pantropical dolphin clusters, where size of circle corresponds to the cluster size, overlaid onto depth data.

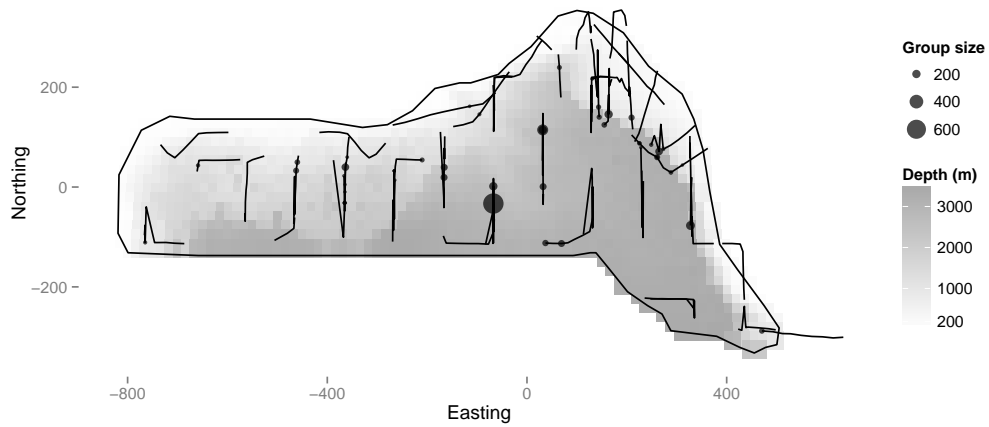


Fig. 3 Predicted abundance of dolphins from the DSM using only depth as an explanatory variable (top) and the model using both depth and location (bottom).

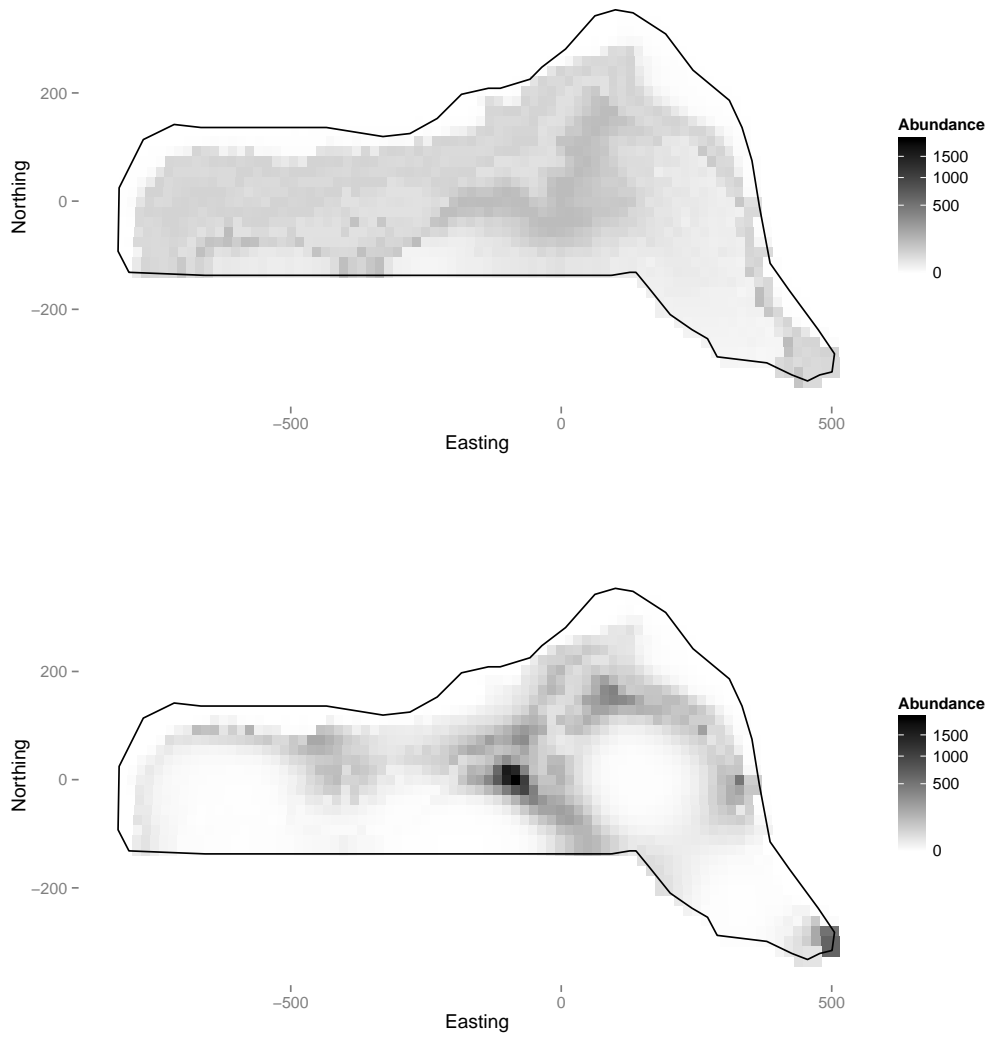


Fig. 4 Plot of the effect on the response of depth, given location (from the model with both depth and location smooths). Note that it is possible to draw a straight line between 750m and 3000m within the confidence band (between the dashed lines), so the wiggles in the smooth may not be indicative of 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.

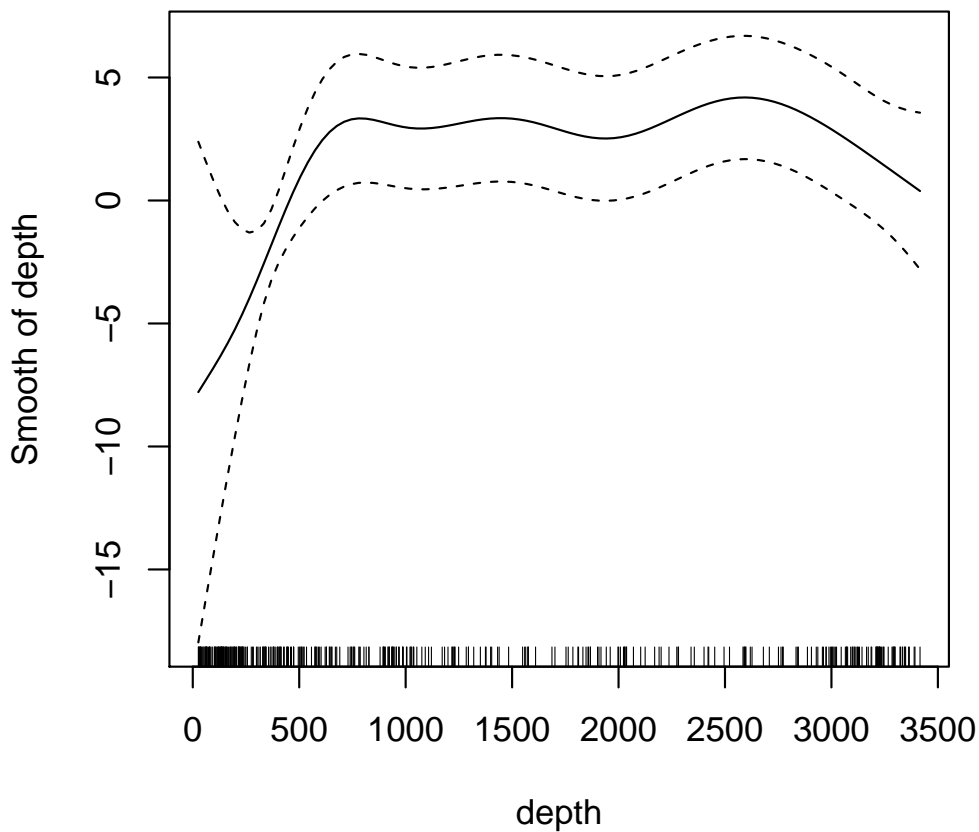


Fig. 5 Map of the coefficients of variation for the model with smooths of both depth and location. Uncertainty was estimated using the variance propagation method of Williams *et al.* (2011). As might be expected, there is high uncertainty where there is low sampling effort (Fig. 2).

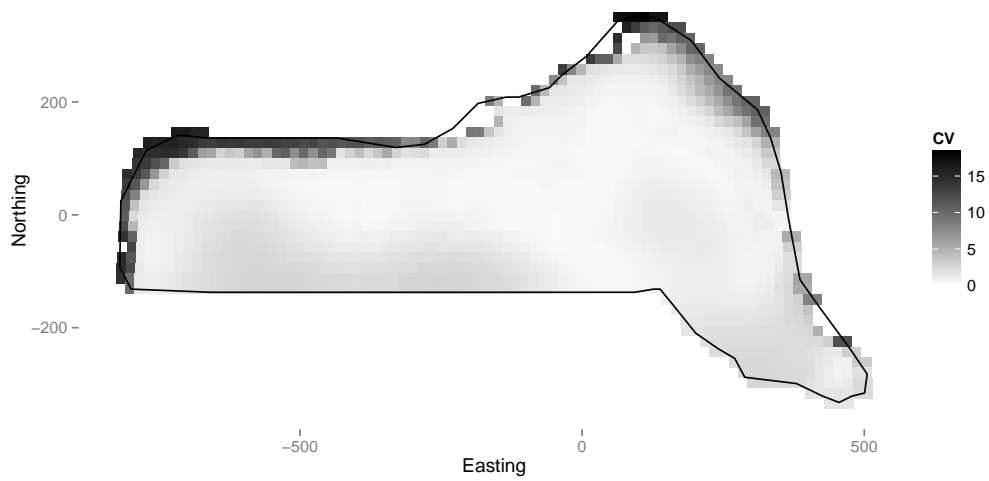


Fig. 6 Flow diagram showing the modelling process for creating a density surface model.

