

1 **Running title:** Spatial models for distance sampling
2 **Number of words:** ~4373
3 **Number of tables:** 0
4 **Number of figures:** 6
5 **Number of references:** 33

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.
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 offer a comparison of recent advances in the field and consider the likely directions of future research. 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 and are largely implemented in the popular Windows package Distance (or are soon to be incorporated).
5. Density surface modelling enables applied ecologists to reliably estimate abundances and create maps of animal/plant distribution. Such models can also be used to investigate the relationships between distribution and environmental covariates.

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 include any
50 model that includes spatially referenced covariates, not just smooths of loc-
51 ation. This article is concerned with combining spatial modelling techniques
52 with distance sampling (Buckland *et al.*, 2001, 2004).

53 Distance sampling takes plot sampling (counting all the individuals or
54 groups of objects within a strip or circle) and extends it to the case where
55 detection is not certain. Observers move along lines or stand at points and
56 record the distance from the line or point to the object of interest (y). These
57 distances are used to estimate the *detection function*, $g(y)$ (for example,
58 Fig. 2), by modelling the decrease in detectability with increasing distance
59 from the line or point (conventional distance sampling, CDS). The detection
60 function may also include covariates (multiple covariate distance sampling,
61 MCDS; Marques *et al.*, 2007) which affect the scale of the detection function.
62 From the fitted detection function, the probability of detection can be estim-

ated. The estimated probability that an animal is detected, \hat{p}_i , can then be used to estimate abundance as

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

where A is the area of the study region, a is the area covered by the survey (i.e., the sum of the areas of all of the strips/circles) and the summation takes place over the n observed individuals (Buckland *et al.*, 2001, Chapter 3). In general, distance sampling is more efficient than plot sampling because a much higher proportion of observations can be used in the analysis. Often up to half the observations in a plot sampling data set are discarded to ensure the assumption of certain detection is met. In contrast, distance sampling uses observations that would have been discarded to model the detection (although typically some detections are discarded beyond a given *truncation distance* during analysis).

When fitting the detection function in a distance sampling analysis, one assumes that the objects of interest are distributed according to some process (Buckland *et al.*, 2001, Section 2.1). It is usually possible to design surveys such that a homogenous process can be assumed so that (with respect to the line/point) objects are distributed uniformly. This can be achieved by ensuring that transects randomly located.

Estimators such as eqn (1) rely on the design of the study to ensure that abundance estimates over the whole study area (scaling up from the covered region) are valid. In contrast this article focusses on *model-based* inference to extrapolate to a larger study area. Specifically, we consider the use of

85 spatially explicit models to investigate the response of biological populations
86 to biotic and abiotic covariates that vary over the study region. A spatially-
87 explicit model can explain the between-transect variation (which is often a
88 large component of the variance in design-based estimates) and so using a
89 model-based approach can lead to smaller variance in estimates of abund-
90 ance. Model-based inference also enables the use of data from opportunistic
91 surveys, for example, incidental data arising from “ecotourism” cruises (Wil-
92 liams *et al.*, 2006).

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

103 In this article we review the current landscape of spatial modelling of
104 distance sampling data, illustrating some recent developments most useful to
105 applied ecologists. The methods discussed have been available in the popular
106 Windows application Distance (Thomas *et al.*, 2010) for some time but the
107 recent advances covered here have been implemented in a new R package,
108 **dsm** (Miller *et al.*, 2013) and are soon to be incorporated into Distance.

109 Throughout this article a motivating data set is used to illustrate the

110 methods. These data are from a combination of several shipboard surveys
 111 conducted on several cetacean species in the Gulf of Mexico. We investigate
 112 47 observations of groups of pantropical spotted dolphins (*Stenella atten-*
 113 *uata*); group size was recorded, as well as the Beaufort sea state at the time
 114 of the observation. Coordinates for each observation and bathymetry data
 115 were available as covariates for the analysis. A complete example analysis is
 116 provided as an online appendix. The data used in the analysis are available
 117 in the `dsm` package and `Distance`.

118 The rest of the article follows this structure: we first introduce the density
 119 surface modelling approach of Hedley & Buckland (2004); explain how to
 120 estimate abundance and uncertainty; describe recent advances and provide
 121 practical advice regarding model fitting, formulation and checking. Before
 122 concluding, we review alternative (but less mature) methods which take a
 123 more direct approach to modelling spatial distance sampling data.

124 Density surface modelling

125 This section focuses on modelling the density/abundance estimation stage
 126 of distance sampling, using the “count model” of Hedley & Buckland (2004),
 127 which we refer to as *density surface modelling* (DSM). Both line and point
 128 transects can be used but if lines are used then they are split into contigu-
 129 ous *segments* (indexed by j), which are of length l_j . Segments should be small
 130 enough such that neither density of objects or covariate values vary appre-
 131 ciably within a segment (usually making the segments approximately square,
 132 $2w \times 2w$, is sufficient). Count or estimated abundance is then modelled as

133 a smooth function of covariates using a generalized additive model (GAM;
 134 e.g. Wood, 2006). For each segment or point, the response is modelled as a
 135 function of environmental covariates that are measured at the segment/point
 136 level (z_{jk} with k indexing the covariates, e.g., location, sea surface temperat-
 137 ure, weather conditions). The area of each segment enters the model as (or
 138 as part of) an offset: the area of segment j is $A_j = 2wl_j$ and for point j is
 139 $A_j = w\pi^2$ (where w is the truncation distance).

140 We begin by describing a formulation where only covariates measured
 141 per-segment (e.g. habitat, beaufort sea state) are included in the detection
 142 function. This simple formulation can be rearranged leading to a framework
 143 for the most general case where covariates included in the detection function
 144 are measured at the observation level (i.e. group size, species).

145 COUNT AS RESPONSE

146 The model for the count per segment is:

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

147 where the f_k s are smooth functions of the covariates and β_0 is an intercept
 148 term. Multiplying the segment area (A_j) by the probability of detection (\hat{p}_j)
 149 gives the *effective area* for segment j . If there are no covariates other than
 150 distance in the detection function then the probability of detection is constant
 151 for all objects observed in the segment (i.e., $\hat{p}_j = \hat{p}, \forall j$). The distribution of
 152 n_j can be modelled as overdispersed Poisson, negative binomial, or Tweedie
 153 distribution (see *Recent developments*, below).

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

As well as simply calculating abundance estimates, relationships between covariates and abundance can be illustrated via plots of marginal smooths. The effect of depth on abundance for the dolphin data can be seen in Fig. 4.

ESTIMATED ABUNDANCE AS RESPONSE

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

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

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

The following model is then fitted:

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

where \hat{N}_j , as with n_j , is assumed to follow an overdispersed Poisson, negative

173 binomial, or Tweedie distribution (see *Recent developments*, below). Note
 174 that the offset is now the area rather than effective area of the segment/point.

175 *DSM with covariates at the observation level*

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

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

183 Density can be modelled rather than abundance by not including offset,
 184 but instead dividing the count (or estimated abundance) by the area of the
 185 segment (and weighting observations by the segment areas). We concentrate
 186 on abundance here, see Hedley & Buckland (2004) for further details on
 187 modelling density.

188 PREDICTION

189 Abundance can be predicted for the each cell in a grid over the region in
 190 question and by summing predicted values over corresponding grid cells.
 191 The areas of the prediction cells must be accounted for in the predictions.
 192 Environmental covariates included in the model must be available at each

193 prediction cell at the required resolution (using prediction grid cells that are
 194 smaller than the resolution of the spatially referenced data have no effect on
 195 abundance/density estimates).

196 VARIANCE ESTIMATION

197 Estimating the variance of abundances calculated using a DSM is not straight-
 198 forward: uncertainty from the estimated parameters of the detection function
 199 must be incorporated into the spatial model. A second consideration is that
 200 in a line transect survey, adjacent segments are likely to be correlated; failure
 201 to account for this spatial autocorrelation will lead to artificially low variance
 202 estimates and hence misleadingly narrow confidence intervals.

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

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

- 208 1. Resample (with replacement) the per-segment residuals, store the val-
 209 ues in \mathbf{r}_b .
- 210 2. Refit the model but with the response set to $\hat{\boldsymbol{\eta}} + \mathbf{r}_b$ (where $\hat{\boldsymbol{\eta}}$ are the
 211 fitted values from the original model).
- 212 3. Take the predicted values for the new model and store them.

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

215 abundance estimate (over the whole region of interest or sub-areas) can then
216 be found by combining the variance estimate from the bootstrap procedure
217 with the variance of the probability of detection from the detection function
218 model (using the delta method which assumes that the two components of
219 the variance are independent; Seber, 1982).

220 The above procedure assumes that there is no correlation in space between
221 segments, if many animals are observed in a particular segment then we might
222 expect there to be high numbers in the adjacent segments. A moving block
223 bootstrap (MBB; Efron & Tibshirani, 1993, Section 8.6) can account for some
224 of this spatial autocorrelation in the variance estimation. The segments are
225 grouped together into overlapping blocks, (so if the block size is 5, block
226 one is segments 1, ..., 5, block two is segments 2, ..., 6, and so on). Then,
227 at step (2) above, resamples are taken of the blocks (contiguous collections
228 of segments) rather than individual segments within the transects. Using
229 blocks should account for some of the autocorrelation between the segments,
230 inflating the variances accordingly. However, because the block size dictates
231 the maximum amount of spatial autocorrelation accounted for, this may not
232 fully account for the autocorrelation. These bootstrap procedures can also be
233 modified to take into account detection function uncertainty by generating
234 new distances from the fitted detection function and then re-calculating the
235 offset by fitting a detection function to the new distances.

236 DSM uncertainty can be visualised via a plot of per-cell coefficient of
237 variation obtained by dividing the standard error for each cell by its predicted
238 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 only segment level covariates are in the DSM. Their procedure is as follows:

1. Fit a density surface model.
2. Re-fit the model with an additional term that characterises the uncertainty in the estimation of the detection function (via the derivatives of the probability of detection, \hat{p}).
3. 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 is not only more computationally efficient but also preferable to the moving block boot-

strap from a technical perspective. A moving block bootstrap does not fully account for spatial autocorrelation because when it reallocates blocks of residuals, it does so without considering the dependence between blocks. This can then lead to wide confidence intervals. The confidence intervals produced via variance propagation are narrower than their bootstrap equivalents, while maintaining good coverage (results of a small simulation study are given in Appendix C).

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

EDGE EFFECTS

Previous work (Ramsay, 2002; Wang & Ranalli, 2007; Wood *et al.*, 2008; Scott-Hayward *et al.*, 2013; Miller & Wood, submitted) has highlighted the need to take care when smoothing over areas with complicated boundaries, e.g., those with rivers, peninsulae or islands. If two parts of the domain (either side of a river or inlet, say) are inappropriately linked by the model (i.e. if the distance between the points is measured as a straight line, rather taking into account obstacles) then the boundary feature can be “smoothed across” leading to incorrect inference. Ensuring that a realistic spatial model has been fitted to the data is essential for valid inference. The soap film smoother of Wood *et al.* (2008) is appealing as the model jointly estimates boundary conditions for a complex study area along with the interior smooth. This can be helpful when uncertainty is estimated via a bootstrap as the model helps avoid large, unrealistic predictions which can plague other

286 smoothers (Bravington & Hedley, 2009).

287 Even if the study area does not have a complicated boundary, edge effects
288 can still be problematic. Miller *et al.* (in prep.) show that global smoothers
289 which have unpenalized plane components tend to cause the fitted surface to
290 increase unrealistically as predictions move further away from the locations
291 of survey effort. They suggest the use of Duchon splines (a generalisation of
292 thin plate regression splines) to alleviate the problem.

293 TWEEDIE DISTRIBUTION

294 The Tweedie distribution offers a flexible alternative to the quasi-Poisson
295 and negative binomial distributions as a response distribution when model-
296 ling count data (Candy, 2004). Through the parameter λ , many common
297 distributions arise; varying λ between 1 (Poisson) and 2 (gamma) leads to
298 a random variable which is a sum of M gamma variables where M is Pois-
299 son distributed (Jørgensen, 1987). The distribution does not change appre-
300 ciably when λ is changed by less than 0.1 therefore, a simple line search
301 over the possible values of λ is usually reasonable. Mark Bravington (pers.
302 comm.) suggested plotting the square root of the absolute value of the re-
303 siduals against fitted values; a “flat” plot (points forming a horizontal line)
304 give an indication of a “good” value for λ . We additionally suggest using the
305 metrics described in the next section for model selection.

306 Practical advice

307 Fig. 6 shows a flow diagram of the modelling process for creating a DSM.
308 The diagram shows which methods are compatible with each other and what
309 the options are for modelling a particular data set.

310 In our experience, it is sensible to obtain a detection function that fits
311 the data as well as possible and only after a satisfactory detection function
312 has been obtained, begin spatial modelling. Model selection for the detection
313 function can be performed using AIC and model checking using goodness-
314 of-fit tests given in (Burnham *et al.*, 2004, Section 11.11). If animals occur
315 in groups rather than individually, bias can be incurred due to the higher
316 visibility of larger groups. It may then be necessary to include size as a
317 covariate in the detection function (see Buckland *et al.*, 2001, Section 4.8.2.4).

318 Smooth terms can be selected using (approximate) p -values, as one would
319 usually for a GAM. An additional useful technique for covariate selection is
320 to use an extra penalty for each term in the GAM allowing smooth terms to
321 be removed from the model during fitting (illustrated in the example ana-
322 lysis; Wood, 2011). Smoothness selection is performed by generalized cross
323 validation (GCV) score, UnBiased Risk Estimator (UBRE) or REstricted
324 Maximum Likelihood (REML) score. When model covariates are effectively
325 functions of one another (e.g. depth could be written as a function of loc-
326 ation) GCV and UBRE can suffer from optimisation failures (Wood, 2006,
327 Section 4.5.3), this can lead to unstable models (Wood, 2011). To avoid these
328 issues REML is recommended for smoothness selection when many spatially-
329 referenced covariates are used. A significant drawback is that REML scores

330 can only be used to compare models with the same fixed effects (i.e. linear
331 terms; Wood, 2011), though the p -value and additional penalty techniques
332 described above can be used to select model terms. We highly recommend
333 the use of standard GAM diagnostic plots; Wood (2006) provides further
334 practical information on GAM model selection and fitting.

335 In the analysis of the dolphin data, we included a smooth of location. This
336 not only nearly doubles the percentage deviance explained (27.3% to 52.7%),
337 it also allows us to account for spatial autocorrelation (in a primitive way).
338 One can see this when comparing the two plots in Fig. 3 and the plot of the
339 depth (Fig. 1), the plot of the model containing only a smooth of depth looks
340 very similar to the raw plot of the depth data. A smooth of an environment-
341 level covariate such as depth can be very useful for assessing the relationships
342 between abundance and the covariate (as in Fig. 4). Caution should be
343 employed when interpreting smooth relationships and abundance estimates,
344 especially if there are gaps over the range of covariate values. Large counts
345 may occur at a high value of depth but if no further observations occur at
346 such a high value, then investigators should be skeptical of any relationship.
347 A smooth of location can be useful although limiting the “wigglyness” of
348 smooths of spatial location (by limiting their basis size) can be a useful
349 way of restricting their influence whilst still allowing them to “mop up” the
350 residual spatial correlation in the data (see the example analysis).

351 In the analysis presented here we have converted spatial location from
352 latitude and longitude to kilometres from the centre of the survey region
353 at $(27.01^\circ, -88.3^\circ)$. This is because the bivariate smoother used (the thin
354 plate spline; Wood, 2003) is isotropic: the wigglyness of the smoother in each

355 direction is treated equally. Moving one degree in latitude is not the same
356 as moving one degree in longitude and so using kilometres from the centre
357 of the study region makes the covariates isotropic (using SI units throughout
358 would also remove the need for conversion).

359 Direct modelling of the spatial point process

360 Rather than use a GAM to model the spatially explicit part of the model,
361 two recent articles (Johnson *et al.*, 2010; Niemi & Fernández, 2010) have
362 used a point process (Cox & Isham, 1980) approach (which was formulated
363 by Hedley & Buckland, 2004). In both cases, the density of the objects is
364 described by an intensity function, which can include spatially-referenced
365 covariates.

366 Johnson *et al.* (2010) proposed a point process-based model for distance
367 sampling data. They first assumed that the locations of all individuals in
368 the survey area (not just those observed) form a realisation of a Poisson
369 process. Parameters of the intensity function were then estimated via stand-
370 ard maximum likelihood methods for point processes (Baddeley & Turner,
371 2000). In contrast to Hedley & Buckland (2004), all parameters were estim-
372 ated jointly so uncertainty from both the spatial pattern and the detection
373 function was incorporated into variance estimates of the abundance. This
374 also ensures that correlations between the detection function and underlying
375 point process are estimated correctly (and do not falsely inflate or deflate
376 variance estimates). The authors also addressed the issue of overdispersion
377 unmodelled by spatial covariates (i.e. counts that do not follow a Poisson

378 mean-variance relationship) using a post-hoc correction factor.

379 Niemi & Fernández (2010) also used Poisson processes but incorporated
380 them into a fully Bayesian approach. Model fitting proceeded in two stages:
381 first the detection function was fitted, then the spatial model (via MCMC)
382 assuming the detection function parameters were known, so detection func-
383 tion uncertainty was not incorporated in the spatial model (an extension that
384 incorporates uncertainty is, however, feasible).

385 Both of the above Poisson process models do not account for group size,
386 but both state that this could be included by considering a marked point
387 process (Cox & Isham, 1980, Section 5.5). Both methods offer direct mod-
388 elling of the point process, although with some drawbacks compared to the
389 methodology of Hedley & Buckland (2004). It should be noted that the loss
390 of efficiency from using DSM is not large (Buckland *et al.*, 2004, p. 313)
391 because distances of detected objects from the line contain little information
392 about spatial variation due to the width of the transects relative to their
393 lengths and how small circles are compared to the study area.

394 A final example of direct modelling of density is given in Royle *et al.*
395 (2004). The authors formulated an unconditional likelihood per-point/line,
396 which is a function of the unobserved transect abundances. These unobserved
397 abundances were treated as (Poisson or negative binomial) random effects,
398 which were then integrated out to give a per-transect likelihood which is a
399 function only of detection function parameters and parameters of the random
400 effects (linear functions of the environmental covariates). Due to the mul-
401 tinomial nature of the per-transect likelihood proposed, distance data must
402 be binned, resulting in a loss of information. Although an arbitrarily large

403 number of bins could be used as an approximation to continuous data, this
404 is potentially computationally intensive.

405 Discussion

406 The use of model-based inference for determining abundance and spatial dis-
407 tribution from distance sampling data presents new opportunities in the field
408 of population assessment. Inference from a sample of sightings to a popula-
409 tion in a study area does not have to depend upon a random sample design,
410 and therefore data collected from "platforms of opportunity" (Williams *et al.*,
411 2006) can be used.

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

419 The field is quickly evolving to allow modelling of more complex data
420 building on the basic ideas of density surface modelling. We expect to see
421 large advances in two areas: temporal inferences and the handling of spa-
422 tial correlation. These should become more mainstream as modern spatio-
423 temporal modelling techniques are adopted. Petersen *et al.* (2011) provided
424 a very basic framework for temporal modelling; their model included "be-
425 fore" and "after" smooth terms to quantify the impact of the construction

426 of an offshore windfarm. Spatial autocorrelation can be accounted for via
427 approaches that explicitly introduce correlations such as generalized estim-
428 ating equations (GEEs; Hardin & Hilbe, 2003) or via mechanisms such as
429 that of Skaug (2006), which allowed observations to cluster according to one
430 of several states (such as high vs low density patches, possibly in response to
431 temporary agglomerations of prey, although the mechanism is unimportant).
432 These advances should assist both modellers and wildlife managers to make
433 optimal conservation decisions.

434 Recent advances in Bayesian computation (INLA; Rue et al, 2009), make
435 one-step, Bayesian, density surface models computationally feasible (as INLA
436 is an alternative to MCMC). We anticipate that such a direct modelling
437 technique will dominate future developments in the field.

438 Density surface modelling allows wildlife managers to make best use of
439 the available spatial data to understand patterns of abundance, and hence
440 make better conservation decisions (e.g., about reserve placement). The re-
441 cent advances mentioned here increase the reliability of the outputs from a
442 modelling exercise, and hence the efficacy of these decisions. Density surface
443 modelling from survey data is an active area of research, and we look forward
444 to further improvements and extensions in the near future.

445 Acknowledgments

446 DLM wishes to thank Mark Bravington and Sharon Hedley for their detailed
447 discussions and for providing code for their variance propagation method.
448 Funding for the implementation of the recent advances into the `dsm` package

449 and Distance software came from the US Navy, Chief of Naval Operations
450 (Code N45), grant number N00244-10-1-0057.

References

- Baddeley, A. & Turner, R. (2000) Practical maximum pseudolikelihood for spatial point patterns. *Australian & New Zealand Journal of Statistics*, **42**, 283–322.
- Bravington, M. & Hedley, S.L. (2009) Antarctic minke whale abundance estimates from the second and third circumpolar IDCR/SOWER surveys using the SPLINTR model. Paper SC/61/IA14, International Whaling Commission Scientific Committee.
- Buckland, S.T., anderson, D.R., Burnham, K.P., Laake, J.L., Borchers, D.L. & Thomas, L. (2001) *Introduction to Distance Sampling*. Oxford University Press.
- Buckland, S.T., anderson, D.R., Burnham, K.P., Laake, J.L., Borchers, D.L. & Thomas, L. (2004) *Advanced Distance Sampling*. Oxford University Press.
- Burnham, K.P., Buckland, S.T., Laake, J.L., Borchers, D.L., Marques, T.A., Bishop, J.R. & Thomas, L. (2004) Further topics in distance sampling. *Advanced Distance Sampling* (eds. S.T. Buckland, D.R. anderson, K.P. Burnham, J.L. Laake, D.L. Borchers & L. Thomas). Oxford University Press.
- Candy, S. (2004) Modelling catch and effort data using generalised linear models, the Tweedie distribution, random vessel effects and random stratum-by-year effects. *Ccamlr Science*, **11**, 59–80.
- Cox, D.R. & Isham, V. (1980) *Point Processes*. Monographs on Applied Probability and Statistics. Chapman and Hall. ISBN 9780412219108.
- Efron, B. & Tibshirani, R.J. (1993) *An Introduction to the Bootstrap*. Chapman & Hall/CRC. ISBN 9780412042317.
- Halpin, P., Read, A., Fujioka, E., Best, B., Donnelly, B., Hazen, L., Kot, C., Urian, K., LaBrecque, E., Dimatteo, A., Cleary, J., Good, C., Crowder, L. & Hyrenbach, K.D. (2009) OBIS-SEAMAP: The World Data Center for Marine Mammal, Sea Bird, and Sea Turtle Distributions. *Oceanography*, **22**, 104–115.
- Hardin, J. & Hilbe, J. (2003) *Generalized Estimating Equations*. Chapman and Hall/CRC, London, UK.
- Hedley, S.L. & Buckland, S.T. (2004) Spatial models for line transect sampling. *Journal of Agricultural, Biological, and Environmental Statistics*, **9**, 181–199.
- Johnson, D.S., Laake, J.L. & Ver Hoef, J.M. (2010) A model-based approach for making ecological inference from distance sampling data. *Biometrics*, **66**, 310–318.

- 484 Jørgensen, B. (1987) Exponential dispersion models. *Journal of the Royal Statist-*
485 *ical Society. Series B, Statistical Methodology*, **49**, 127–162.
- 486 Marques, T.A., Thomas, L., Fancy, S. & Buckland, S.T. (2007) Improving estimates
487 of bird density using multiple-covariate distance sampling. *The Auk*, **124**, 1229–
488 1243.
- 489 Miller, D.L., Rexstad, E., Burt, L., Bravington, M.V. & Hedley., S. (2013) *dsm:*
490 *Density surface modelling of distance sampling data*. R package version 2.0.1.
491 URL <http://cran.r-project.org/package=Distance>
- 492 Miller, D.L., Jones, E. & Matthiopoulos, J. (in prep.) Reliable spatial smoothing
493 without edge effects. pp. 1–8.
- 494 Miller, D.L. & Wood, S.N. (submitted) Finite area smoothing with generalized
495 distance splines. pp. 1–27.
- 496 Niemi, A. & Fernández, C. (2010) Bayesian Spatial Point Process Modeling of Line
497 Transect Data. *Journal of Agricultural, Biological, and Environmental Statistics*,
498 **15**, 327–345.
- 499 Petersen, I.K., MacKenzie, M.L., Rexstad, E.A., Wisz, M.S. & Fox, A.D. (2011)
500 Comparing pre- and post-construction distributions of long-tailed ducks *Clangula*
501 *hyemalis* in and around the Nysted offshore wind farm, Denmark: a quasi-
502 designed experiment accounting for imperfect detection, local surface features
503 and autocorrelation. CREEM technical report, University of St Andrews.
- 504 Ramsay, T. (2002) Spline smoothing over difficult regions. *Journal of the Royal*
505 *Statistical Society. Series B, Statistical Methodology*, **64**, 307–319.
- 506 Royle, J., Dawson, D. & Bates, S. (2004) Modeling abundance effects in distance
507 sampling. *Ecology*, **85**, 1591–1597.
- 508 Rue, H., Martino, S. & Chopin, N. (2009) Approximate Bayesian inference for
509 latent Gaussian models by using integrated nested Laplace approximations. *J.*
510 *R. Statist. Soc. B*, **71**, 319–392.
- 511 Scott-Hayward, L.A.S., MacKenzie, M.L., Donovan, C.R., Walker, C.G. & Ashe, E.
512 (2013) Complex Region Spatial Smoother (CReSS). *Journal of Computational*
513 *and Graphical Statistics*.
- 514 Seber, G.A.F. (1982) *The Estimation of Animal Abundance and Related Paramet-*
515 *ers*. ISBN 9781930665552.
- 516 Skaug, H.J. (2006) Markov modulated Poisson processes for clustered line transect
517 data. *Environmental and Ecological Statistics*, **13**, 199–211.

- 518 Thomas, L., Buckland, S.T., Rexstad, E.A., Laake, J.L., Strindberg, S., Hedley,
519 S.L., Bishop, J.R., Marques, T.A. & Burnham, K.P. (2010) Distance software:
520 design and analysis of distance sampling surveys for estimating population size.
521 *Journal of Applied Ecology*, **47**, 5–14.
- 522 Wang, H. & Ranalli, M. (2007) Low-rank smoothing splines on complicated do-
523 mains. *Biometrics*, **63**, 209–217.
- 524 Williams, R., Hedley, S.L., Branch, T.A., Bravington, M.V., Zerbini, A.N. & Find-
525 lay, K.P. (2011) Chilean blue whales as a case study to illustrate methods to
526 estimate abundance and evaluate conservation status of rare species. *Conserva-
527 tion Biology*, **25**, 526–535.
- 528 Williams, R., Hedley, S.L. & Hammond, P. (2006) Modeling distribution and
529 abundance of Antarctic baleen whales using ships of opportunity. *Ecology and
530 Society*, **11**, 1.
- 531 Wood, S.N. (2003) Thin plate regression splines. *Journal of the Royal Statistical
532 Society. Series B, Statistical Methodology*, **65**, 95–114.
- 533 Wood, S.N. (2006) *Generalized Additive Models: An introduction with R*. Chapman
534 & Hall/CRC.
- 535 Wood, S.N. (2011) Fast stable restricted maximum likelihood and marginal like-
536 lihood estimation of semiparametric generalized linear models. *Journal of the
537 Royal Statistical Society. Series B, Statistical Methodology*, **73**, 3–36.
- 538 Wood, S.N., Bravington, M.V. & Hedley, S.L. (2008) Soap film smoothing. *Journal
539 of the Royal Statistical Society. Series B, Statistical Methodology*, **70**, 931–955.

Figures

Fig. 1 The region, transect centrelines and location of detected pantropical dolphin groups, where size of circle corresponds to the group size, overlaid onto depth data.

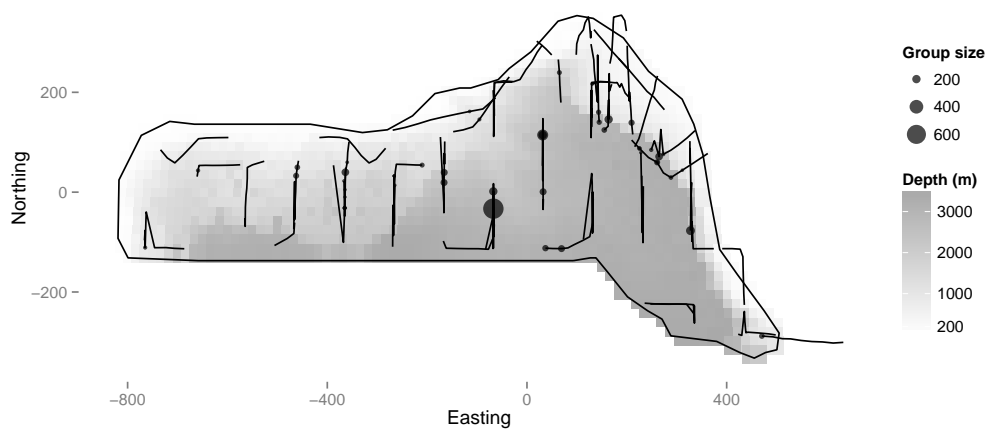


Fig. 2 Estimated detection function for pantropical dolphin groups overlaid onto the scaled histogram of observed distances. Distances are recorded in metres.

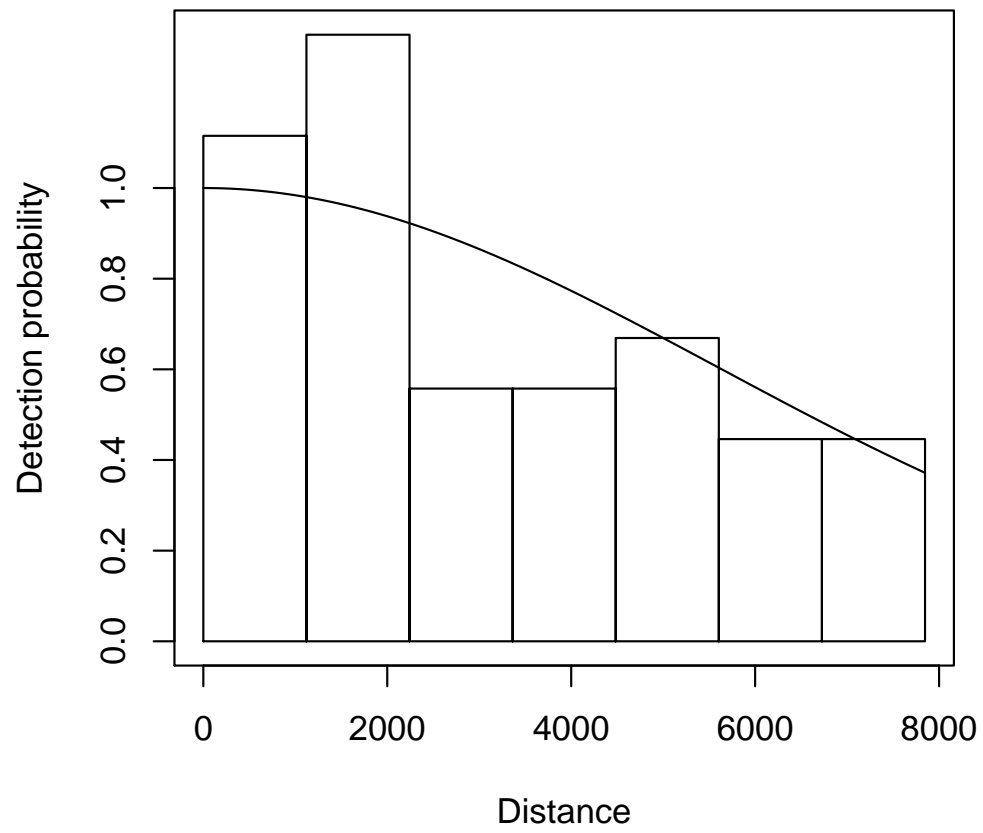


Fig. 3 Predicted abundance of dolphins from the model using only depth as an explanatory variable (top) and the model using both depth and location (bottom).

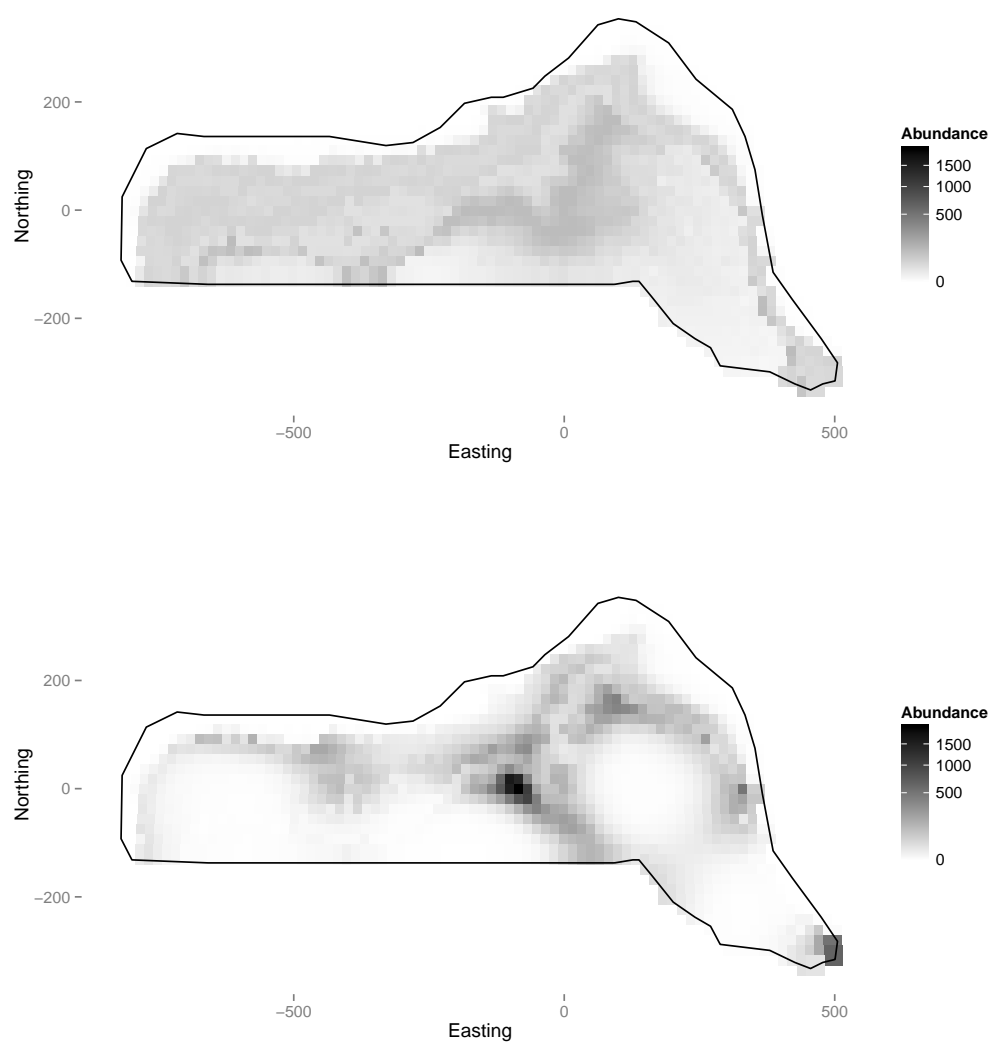


Fig. 4 Plot of the effect on the response of depth (from the model with both depth and location smooths), note that it is possible to draw a straight line between 750m and 3000m within the confidence band (between the dashed lines), so the wiggles in the smooth may not be indicative of any relationship.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

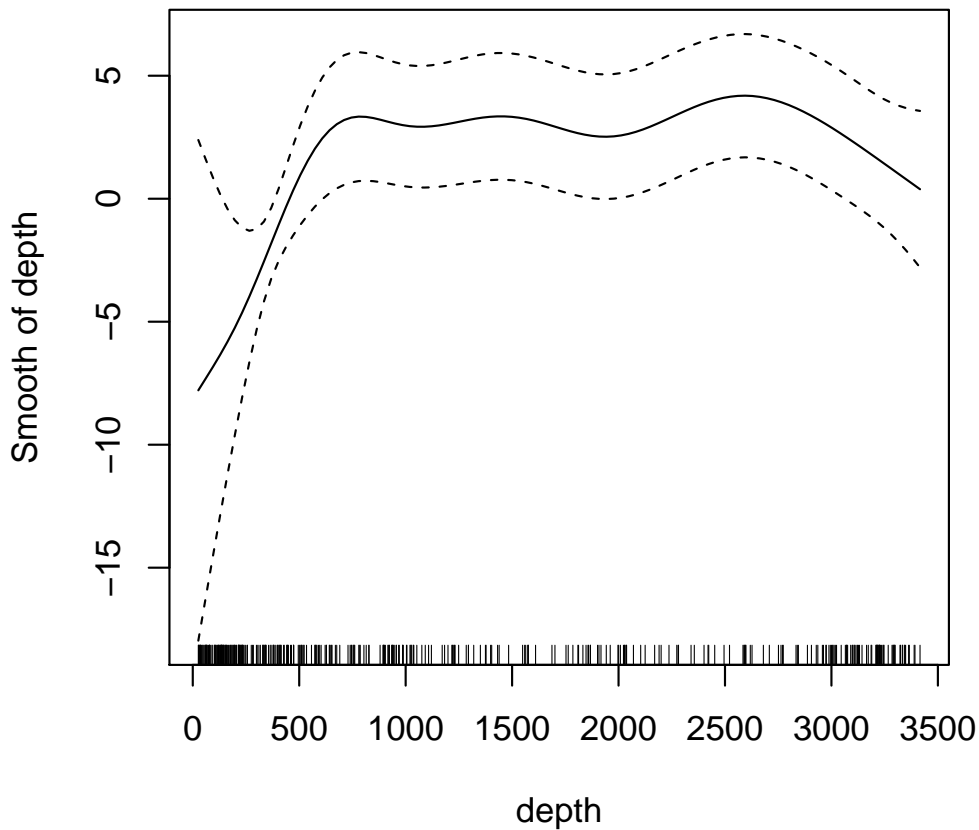


Fig. 5 Map of the coefficients of variation for the model with smooths of both depth and location. Uncertainty was estimated using the variance propagation method of Williams *et al.* (2011). As might be expected, there is high uncertainty where there is low sampling effort (Fig. 1).

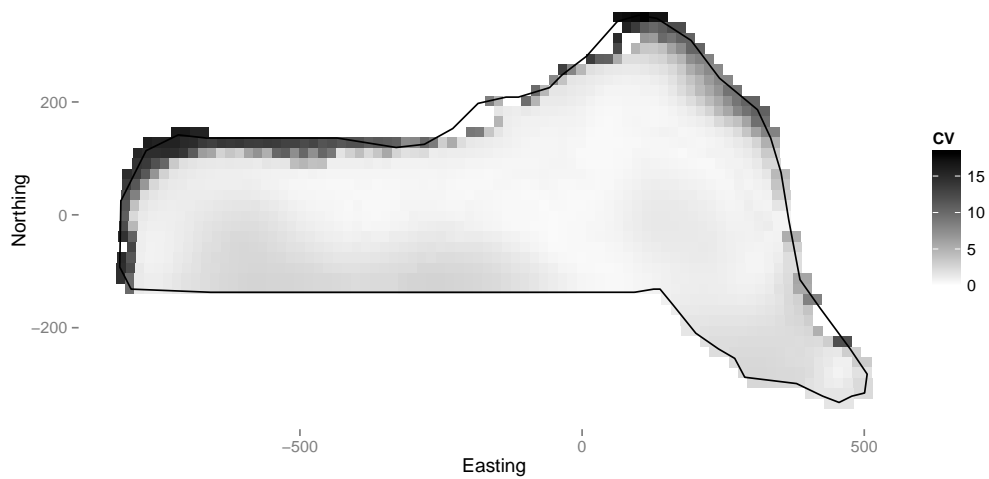


Fig. 6 Flow diagram showing the modelling process for creating a density surface model.

