# Spatial models for distance sampling data: recent developments and future directions

8 **David L. Miller**[1*],      **M. Louise Burt**[2],
9       **Eric A. Rexstad**[2],       **Len Thomas**[2].

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.

1

# Introduction

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

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

2

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

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

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 the observations that would have been discarded to model the detection (typically data 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) 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. By contract this article focusses on *model-based* inference to extrapolate to a larger study area. Specifically, we consider the

use spatially explicit models to investigate the response of biological populations to biotic and abiotic covariates that vary over the study region. A spatially-explicit model can explain the between transect variation (which is often a large component of the variance in design-based estimates) and so using a model-based approach can lead to smaller variance in estimates of abundance. Model-based inference also enables the use of data from opportunistic surveys, for example, incidental data arising from "ecotourism" cruises (Williams *et al.*, 2006).

Our aims in creating a spatial model of a biological population are usually two-fold: (i) estimating overall abundance and (ii) investigating the relationship between abundance and environmental covariates. As with any predictions that are outside the range of the data, one should heed the usual warnings regarding extrapolation. For example, if a model contains elevation as a covariate, predictions at high, unsampled elevations are unlikely to be reliable. Frequently, maps of abundance or density are required and any spurious predictions can be visually assessed, as well as by plotting a histogram of the predicted values. A sensible definition of the region of interest avoids prediction outside the range of the data.

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

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

4

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

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

# Density surface modelling

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

<sub>133</sub> a smooth function of covariates using a generalized additive model (GAM;

<sub>134</sub> e.g. Wood, 2006). For each segment or point, the response is modelled as a

<sub>135</sub> function of environmental covariates that are measured at the segment/point

<sub>136</sub> level ($z_{jk}$ with $k$ indexing the covariates, e.g., location, sea surface temperat-

<sub>137</sub> ure, weather conditions). The area of each segment enters the model as (or

<sub>138</sub> as part of) an offset: the area of segment $j$ is $A_j = 2wl_j$ and at point $j$ is

<sub>139</sub> $A_j = w\pi^2$ (where $w$ is the truncation distance).

<sub>140</sub> COUNT AS RESPONSE

<sub>141</sub> The model for the count per segment is:

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

<sub>142</sub> where the $f_k$s are smooth functions of the covariates and $\beta_0$ is an intercept

<sub>143</sub> term. Multiplying the segment area $(A_j)$ by the probability of detection $(\hat{p}_j)$

<sub>144</sub> gives the *effective area* for segment $j$. If there are no covariates other than

<sub>145</sub> distance in the detection function then the probability of detection is constant

<sub>146</sub> for all objects observed in the segment (i.e., $\hat{p}_j = \hat{p}, \forall j$). The distribution of

<sub>147</sub> $n_j$ can be modelled as overdispersed Poisson, negative binomial, or Tweedie

<sub>148</sub> distribution (see *Recent developments*, below).

<sub>149</sub> Fig. 1 shows the raw observations of the dolphin data, along with the

<sub>150</sub> transect lines, overlaid on the depth data. A half-normal detection function

<sub>151</sub> was fitted to the distances and is shown in Fig. 2. Fig. 3 shows a DSM fitted

<sub>152</sub> to the dolphin data. The top panel shows predictions from a model where

<sub>153</sub> depth was the only covariate, the bottom panel shows predictions where a

<sub>6</sub>

<sup>154</sup> (bivariate) smooth of spatial location was also included. The latter had a
<sup>155</sup> considerably lower GCV score (39.12 vs 48.46) so would be selected as our
<sup>156</sup> "best" model.

<sup>157</sup> As well as simply calculating abundance estimates, relationships between
<sup>158</sup> covariates and abundance can be illustrated via plots of marginal smooths.
<sup>159</sup> The effect of depth on abundance for the dolphin data can be seen in Fig. 4.

<sup>160</sup> ESTIMATED ABUNDANCE AS RESPONSE

<sup>161</sup> An alternative to modelling counts is to use the per-segment/circle abund-
<sup>162</sup> ance using distance sampling estimates as the response. In this case we
<sup>163</sup> replace $n_j$ by:

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

<sup>164</sup> where $R_j$ is the number observations in segment $j$ and $s_{jr}$ is the size of the
<sup>165</sup> $r^{\text{th}}$ group in segment $j$ (if the animals occur individually then $s_{jr} = 1, \forall j, r$).

<sup>166</sup> The following model is then fitted:

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

<sup>167</sup> where $\hat{N}_j$, as with $n_j$, is assumed to follow an overdispersed Poisson, negative
<sup>168</sup> binomial, or Tweedie distribution (see *Recent developments*, below). Note
<sup>169</sup> that the offset is now the area rather than effective area of the segment/point.

7

*DSM with covariates at the observation level*

171 The above models consider the case where the covariates are measured at

172 the segment/point level. Often covariates ($z_{ij}$, for individual/group $i$ and

173 segment/point $j$) are collected on the level of observations; for example sex

174 or group size of the observed object or identity of the observer. In this case

175 the probability of detection is a function of the object (individual or group)

176 level covariates $\hat{p}(z_i)$. Object level covariates can be incorporated into the

177 model by adopting the following estimator of the per-segment abundance:

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

178 By not including an offset, but instead dividing the count (or estimated

179 abundance) by the area of the segment, we can also model density rather

180 than abundance. We concentrate on abundance here, see Hedley & Buckland

181 (2004) for further details on modelling density.


182 PREDICTION

183 Abundance can be predicted for the each cell in a grid over the region in

184 question and by summing predicted values over corresponding grid cells.

185 The areas of the prediction cells must be accounted for in the predictions.

186 Environmental covariates included in the model must be available at each

187 prediction cell at the required resolution (using prediction grid cells that are

188 smaller than the resolution of the spatially referenced data have no effect on

189 abundance/density estimates).

VARIANCE ESTIMATION

191 Estimating the variance of abundances calculated using a DSM is not straight-

192 forward: uncertainty from the estimated parameters of the detection function

193 must be incorporated into the spatial model. A second consideration is that

194 in a line transect survey, adjacent segments are likely to be correlated; failure

195 to account for this spatial autocorrelation will lead to artificially low variance

196 estimates and hence misleadingly narrow confidence intervals.

197 Hedley & Buckland (2004) describe a method of calculating the variance

198 in the abundance estimates using a parametric bootstrap, resampling from

199 the residuals of the fitted model. The bootstrap procedure is as follows.

200 Denote the fitted values for the model to be $\hat{\boldsymbol{\eta}}$. For $b = 1, \ldots, B$ (where

201 $B$ is the number of resamples required).

202 1. Resample (with replacement) the per-segment residuals, store the val-
203 ues in $\mathbf{r}_b$.

204 2. Refit the model but with the response set to $\hat{\boldsymbol{\eta}} + \mathbf{r}_b$ (where $\hat{\boldsymbol{\eta}}$ are the
205 fitted values from the orginal model).

206 3. Take the predicted values for the new model and store them.

207 From the predicted values stored in the last step the variance originating in

208 the spatial part of the model can be calculated. The total variance of the

209 abundance estimate (over the whole region of interest or sub-areas) can then

210 be found by combining the variance estimate from the bootstrap procedure

211 with the variance of the probability of detection from the detection function

212 model (using the delta method which assumes that the two components of

9

the variance are independent; Seber, 1982).

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

DSM uncertainty can be visualised via a plot of per-cell coefficient of variation obtained by dividing the standard error for each cell by its predicted abundance.

10

# 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 as 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

Recent 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 (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 smoothers

12

280 (Bravington & Hedley, 2009).

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

287     Tweedie distribution

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

13

# Practical advice

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

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

Smooth terms can be selected using (approximate) $p$-values in GAM. A useful technique for covariate selection is to use an additional penalty for each term in the GAM allowing smooth terms to be removed from the model during fitting (illustrated in the example analysis; Wood, 2011). Smoothness selection is performed by generalized cross validation (GCV) score, UnBiased Risk Estimator (UBRE) or REstricted Maximum Likelihood (REML) score. When model covariates are effectively functions of one another (e.g. depth could be written as a function of location) GCV and UBRE can suffer from concurvity issues which lead to failures in optimisation (Wood, 2006, Section 4.5.3). The minima in GCV/UBRE tend to have less pronounced minima than REML so an optimal degree of smoothing may not be found, this can lead to unstable models (slight changes in smoothness lead to vastly differ-

14

ent results; Wood, 2011). To avoid these issues REML is recommended for smoothness selection, when many spatially-referenced covariates are used. A significant drawback is that REML scores can only be used to compare models with the same fixed effects (i.e. linear terms; Wood, 2011). We highly recommend the use of standard GAM diagnostic plots. Wood (2006) provides further practical information on GAM model selection and fitting.

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

In the analysis presented we have converted from latitude and longitude to kilometres from the centre of the survey region (27.01, -88.3) because the bivariate smoother used (the thin plate spline; Wood, 2003) is isotropic the

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

# Direct modelling of the spatial point process

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

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

16

relationship) using a post-hoc correction factor.

Niemi & Fernández (2010) also used Poisson processes but incorporate them into a fully Bayesian approach. Model fitting proceeds in two stages: first the detection function is fitted, then the spatial model (via MCMC) assuming the detection function parameters are known, so detection function uncertainty is not incorporated in the spatial model (an extension that incorporates uncertainty is, however, feasible).

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

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

large number of bins could be used as an approximation, this is potentially computationally intensive.

# Discussion

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

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

The field is quickly evolving to allow modelling of more complex data building on the basic ideas of density surface modelling. We expect to see large advances in two areas: temporal inferences and the handling of spatial correlation. These should become more mainstream as modern spatio-temporal modelling techniques are adopted. Petersen *et al.* (2011) provided a very basic framework for temporal modelling; their model included smooth terms both before and after the construction of an offshore windfarm. Spatial autocorrelation can be accounted for via approaches that explicitly intro-

duce correlations such as generalized estimating equations (GEEs; Hardin & Hilbe, 2003) or via mechanisms such as that of Skaug (2006), which allowed observations to cluster according to one of several states (e.g. "feeding" or "transit") taking into account short-term agglomerations ("hot spots"). These advances should assist both modellers and wildlife managers to make optimal conservation decisions.

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

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

# Acknowledgments

# 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 Commisson 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.

Jørgensen, B. (1987) Exponential dispersion models. *Journal of the Royal Statistical Society. Series B, Statistical Methodology*, **49**, 127–162.

Marques, T.A., Thomas, L., Fancy, S. & Buckland, S.T. (2007) Improving estimates of bird density using multiple-covariate distance sampling. *The Auk*, **124**, 1229–1243.

Miller, D.L., Rexstad, E., Burt, L., Bravington, M.V. & Hedley., S. (2013) *dsm: Density surface modelling of distance sampling data*. R package version 2.0.1. URL http://cran.r-project.org/package=Distance

Miller, D.L., Jones, E. & Matthiopoulos, J. (in prep.) Reliable spatial smoothing without edge effects. pp. 1–8.

Miller, D.L. & Wood, S.N. (submitted) Finite area smoothing with generalized distance splines. pp. 1–27.

Niemi, A. & Fernández, C. (2010) Bayesian Spatial Point Process Modeling of Line Transect Data. *Journal of Agricultural, Biological, and Environmental Statistics*, **15**, 327–345.

Petersen, I.K., MacKenzie, M.L., Rexstad, E.A., Wisz, M.S. & Fox, A.D. (2011) Comparing pre- and post-construction distributions of long-tailed ducks *Clangula hyemalis* in and around the Nysted offshore wind farm, Denmark: a quasi-designed experiment accounting for imperfect detection, local surface features and autocorrelation. CREEM technical report, University of St Andrews.

Ramsay, T. (2002) Spline smoothing over difficult regions. *Journal of the Royal Statistical Society. Series B, Statistical Methodology*, **64**, 307–319.

Royle, J., Dawson, D. & Bates, S. (2004) Modeling abundance effects in distance sampling. *Ecology*, **85**, 1591–1597.

Rue, H., Martino, S. & Chopin, N. (2009) Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. *J. R. Statist. Soc. B*, **71**, 319–392.

Scott-Hayward, L.A.S., MacKenzie, M.L., Donovan, C.R., Walker, C.G. & Ashe, E. (2013) Complex Region Spatial Smoother (CReSS). *Journal of Computational and Graphical Statistics*.

Seber, G.A.F. (1982) *The Estimation of Animal Abundance and Related Parameters*. ISBN 9781930665552.

Skaug, H.J. (2006) Markov modulated Poisson processes for clustered line transect data. *Environmental and Ecological Statistics*, **13**, 199–211.

Thomas, L., Buckland, S.T., Rexstad, E.A., Laake, J.L., Strindberg, S., Hedley, S.L., Bishop, J.R., Marques, T.A. & Burnham, K.P. (2010) Distance software: design and analysis of distance sampling surveys for estimating population size. *Journal of Applied Ecology*, **47**, 5–14.

Wang, H. & Ranalli, M. (2007) Low-rank smoothing splines on complicated domains. *Biometrics*, **63**, 209–217.

Williams, R., Hedley, S.L., Branch, T.A., Bravington, M.V., Zerbini, A.N. & Findlay, K.P. (2011) Chilean blue whales as a case study to illustrate methods to estimate abundance and evaluate conservation status of rare species. *Conservation Biology*, **25**, 526–535.

Williams, R., Hedley, S.L. & Hammond, P. (2006) Modeling distribution and abundance of Antarctic baleen whales using ships of opportunity. *Ecology and Society*, **11**, 1.

Wood, S.N. (2003) Thin plate regression splines. *Journal of the Royal Statistical Society. Series B, Statistical Methodology*, **65**, 95–114.

Wood, S.N. (2006) *Generalized Additive Models: An introduction with R* . Chapman & Hall/CRC.

Wood, S.N. (2011) Fast stable restricted maximum likelihood and marginal likelihood estimation of semiparametric generalized linear models. *Journal of the Royal Statistical Society. Series B, Statistical Methodology*, **73**, 3–36.

Wood, S.N., Bravington, M.V. & Hedley, S.L. (2008) Soap film smoothing. *Journal of the Royal Statistical Society. Series B, Statistical Methodology*, **70**, 931–955.

# Figures

**Fig. 1** The survey area for the example dolphin analysis, transect centrelines and observations with size of circle corresponding to the group size overlaid onto depth data.
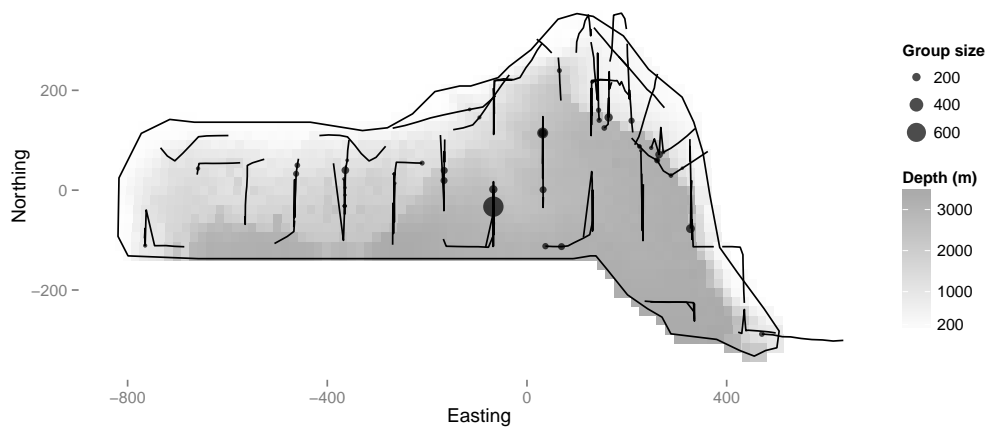
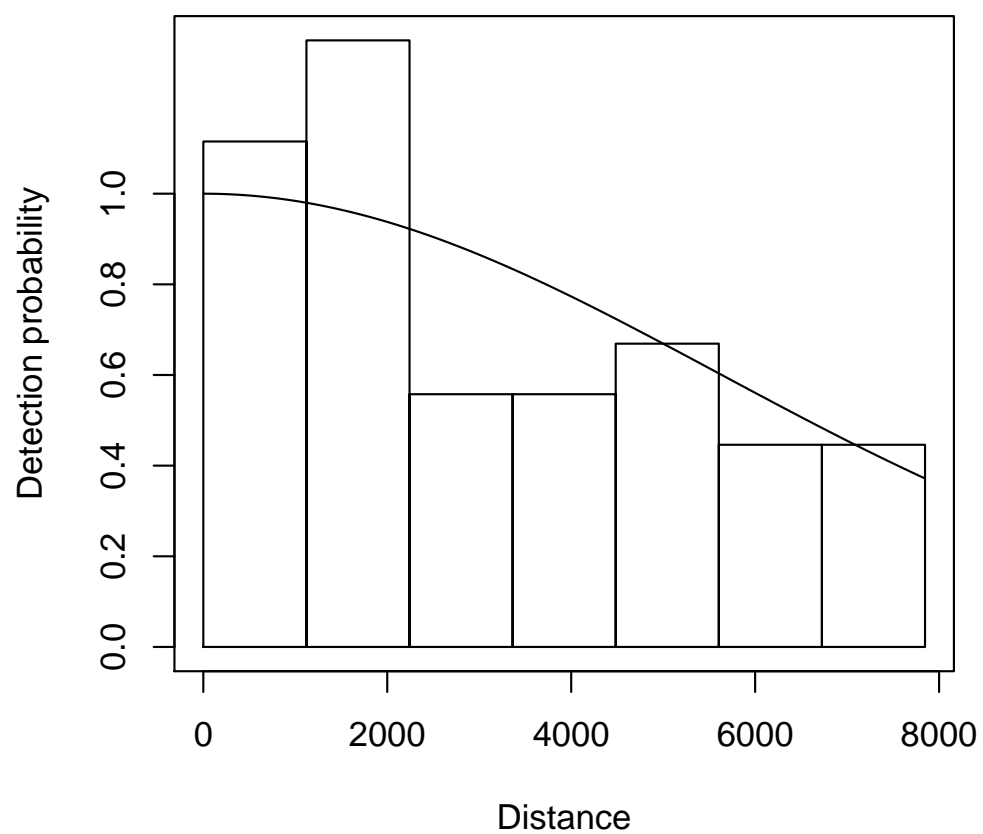**Fig. 2** Histogram of observed distances with detection function fitted to the dolphin data.
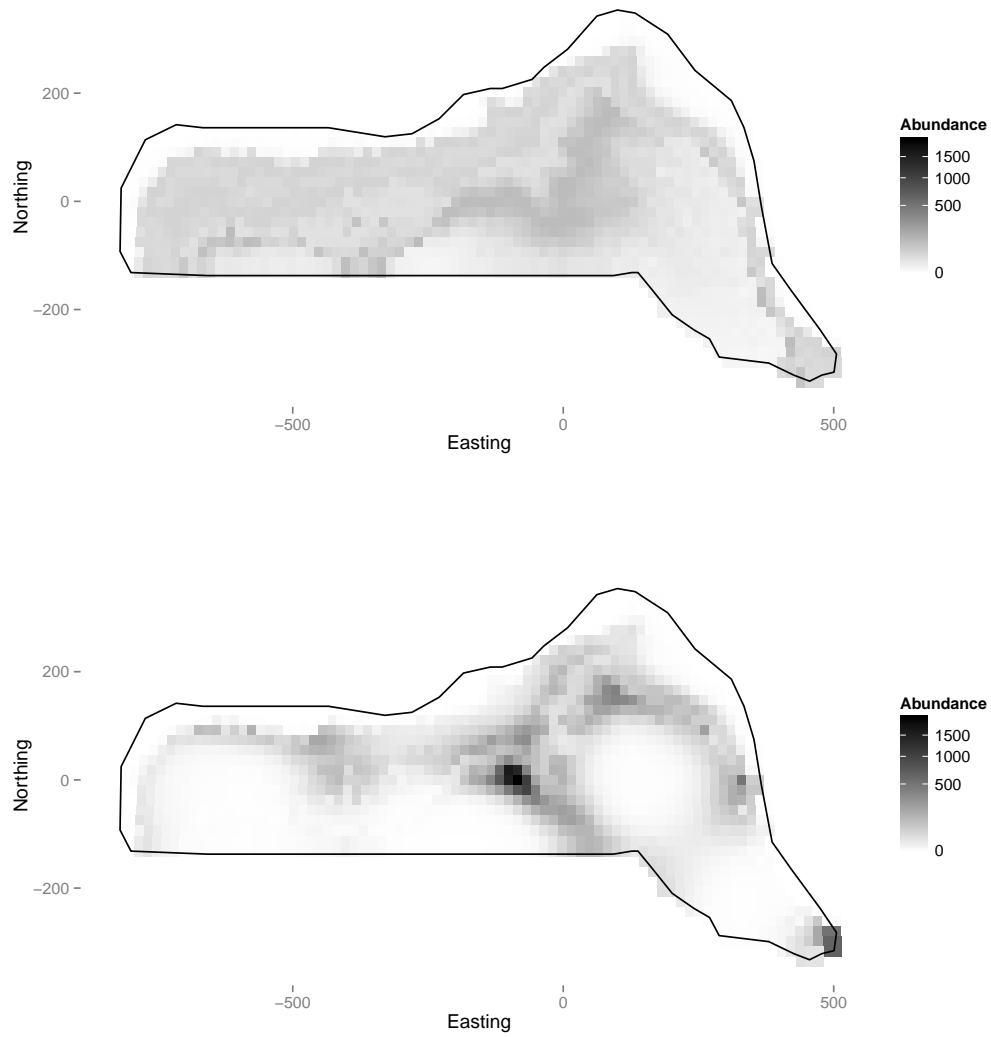
**Fig. 4**   Plot of the effect on the response of depth (from the model with both depth and location smooths), note that it is possible to draw a straight line between 750m and 3000m within the confidence band (between the dashed lines), so the wiggles in the smooth may not be indicative of any relationship. 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.
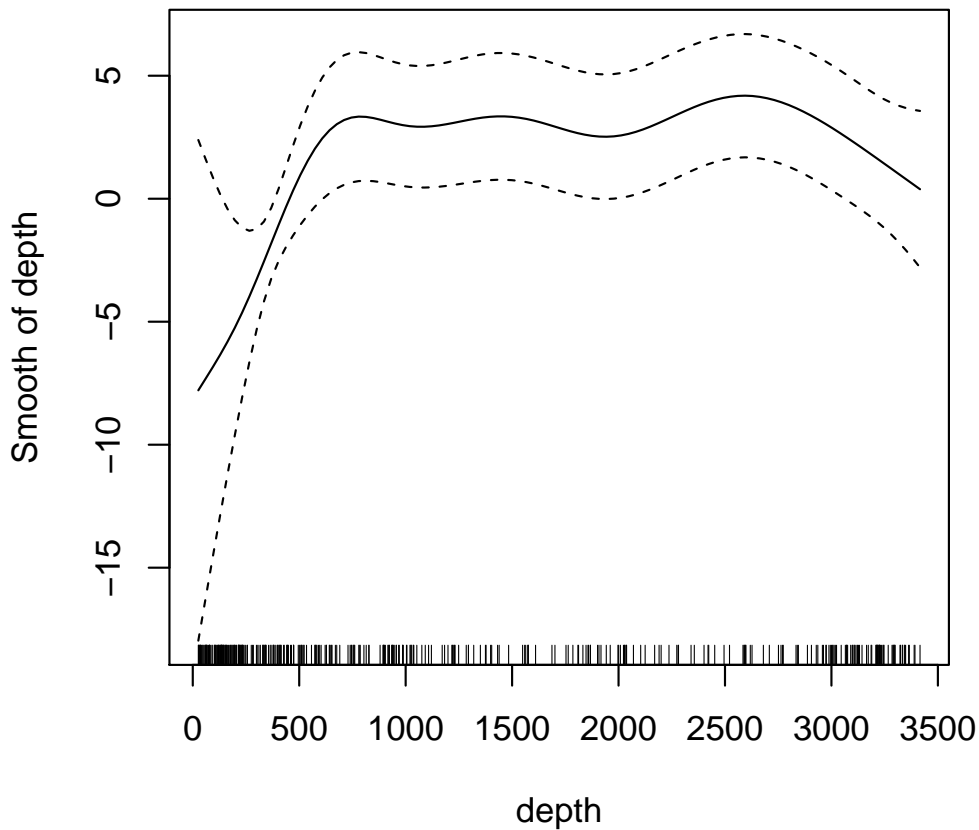
**Fig. 5** Plot of coefficient of variation map for the model with smooths of both depth and location. Uncertainty was estimated using the variance propagation method of Williams *et al.* (2011). As might be expected, there is high uncertainty where there is low sampling effort (Fig. 1).
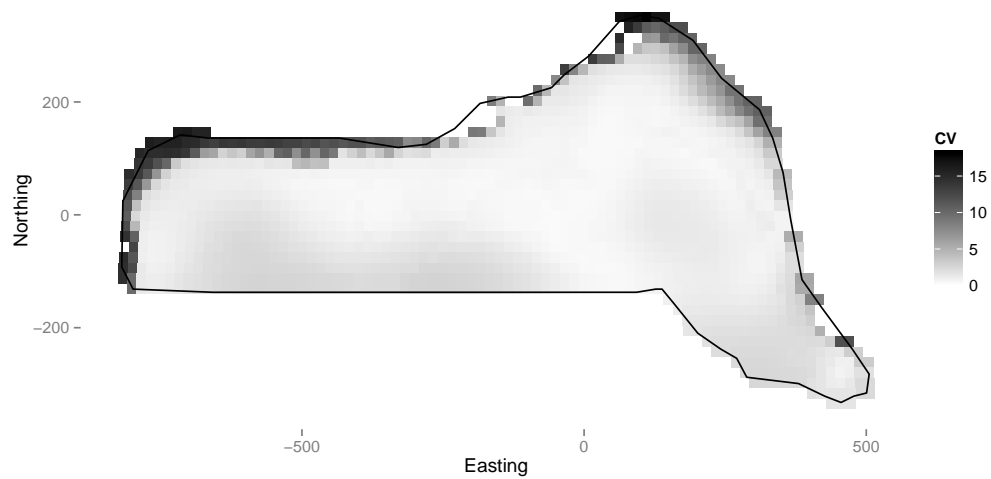
**Fig. 6** Flow diagram showing the modelling process for creating a density surface model.