

Biomathematics and Statistics Scotland

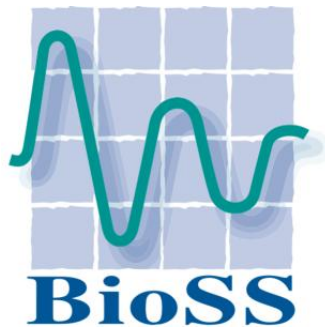
**To carry out tern modelling under the Framework
Agreement C10-0206-0387**

CONTRACT No: C10-0206-0387

Report submitted to:

Joint Nature Conservation Committee

March 2012



Authors:

Mark J Brewer, Jackie M Potts, Elizabeth I Duff,

David A Elston

CONTENTS	PAGE
1. Non Technical Summary	3
2. Introduction	4
3. Data	5
4. Methodology	8
5. Results	14
6. Discussion	19
7. Maps Appendix	20

In addition to this report, there are two further documents associated with this project:

(i) BioSS Terns Report – Results Appendix;

(ii) BioSS Terns Report – Software;

and also ancillary files:

(iii) Spreadsheet files of grid predictions for each of the twelve species/colony combinations;

(iv) R code files for: data exploration; model fitting; grid predictions

(v) cleaned and standardised versions of the original data files.

1. Non-Technical Summary

The Joint Nature Conservation Committee (JNCC) is working on the identification of important marine areas around the UK that are used by five species of tern during the breeding season. For the four larger tern species (Arctic, common, roseate and Sandwich terns), data are available from boat surveys, using both visual tracking and transect survey methods.

Following a competitive tendering process, in December 2011 BioSS was tasked with analysis of the visual tracking data for the four larger species of tern.

The analysis was to use the observed tern data in combination with data simulated from the whole of the sample area, with statistical models attempting to relate the locations of terns with background environmental variables. From these statistical models, predictions were to be made on a grid of data and then mapped, showing areas which are preferred by terns.

We found during analysis that the existing statistical methodology was inadequate for analysis, and that a procedure which accounts for the fact that birds are tracked for different lengths of time was needed.

The environmental data contained sets of variables which were very highly related to each other, and so the analyses can only assess associations, rather than identify drivers. We found that for most species, the most important predictor was the distance to the colony itself, but above and beyond that, other variables were also important. For both the Arctic and common terns, the average sea surface temperature in spring proved an important predictor. With only one colony of roseate terns, we could not draw conclusions of consistency across colonies, but chlorophyll may play a part in addition to distance to the colony and sea surface temperature. For Sandwich terns, the clearest predictor was distance to shore, and while other variables seemed important for particular species, there seemed to be nothing consistently so.

2. Introduction

The problem at hand is thus: given the locations of terns, recorded by tracking (as described in Section 3), how can we learn about the preference of the birds in respect to how they select where to forage (based on the explanatory environmental variables we have), and how can this information be used to make predictions about locations frequented by different tern species?

As we are not sure as to the kind (shape) of relationships between tern preference and the environmental variables – that is, whether they will be linear or non-linear – we will need to consider a class of models which allow flexible regression relationships (Generalised Additive Models, GAMs) as well as those which are less flexible (Generalised Linear Models, GLMs).

The first difficulty is in the form of the data. Since birds fly around and since we wish to distinguish between foraging and commuting behaviour, we cannot simply take a random sample of locations and assess whether these locations were visited by terns, which would be a straightforward logistic regression analysis. Instead, the data represent, in essence, “presence-only” data. Ecologists, taking a lead from Boyce and McDonald (1999), who in turn borrowed from the medical statistics literature, started using so-called case-control designs; in medicine, these studies compare individuals having a condition (the “cases”) with “control” individuals who do not have the condition. In spatial ecology, this equates to the observed presences being the cases, but unless absence can be determined reliably (as might be possible, for example, if we were studying nesting sites), we lack controls. A partial solution has been to sample from the entire possible range for the population at hand, and to use these as controls; this represents a so-called use-availability design. These controls represent pseudo-absences rather than true absences and because the number of controls is determined arbitrarily, we can only estimate relative preference and not absolute preference.

The second problem is that the data consist of repeated observations for a single bird, and that the number of repeats can vary enormously, depending on the length of the trip each bird made and how long the observers were able to track it. If not accounted for properly, this can bias the results. However, thinning the data set to a single observation per individual is not sensible, given the amount of data loss involved and the fact that we are interested in location-specific covariates, although a certain amount of thinning (equally applied to all individuals) was necessary to make the data set manageable for analysis. Previous work in the area (e.g. Aarts *et al.*, 2008) applied “mixed models” to account for between-individual variation, but does not seem to have fully addressed the problem of serial correlation within a track (see later). The tern data set differs from the one studied by Aarts *et al.* (2008) in that we have only one track per individual, so we cannot distinguish variability between individuals from variability between trips. Thus rather than simply applying existing methodology, there was also to be an element of method development.

Another difficulty relates to the spatial nature of the data. Technically there is a temporal element too, but this only applies within a single track, so we consider the spatial aspect. It is well understood (Beale *et al.*, 2010) that failing to deal with spatial autocorrelation – that is, the fact that data from locations close in space are typically more similar than data further apart – can lead to incorrect statistical inferences; in particular, significant associations can be claimed which are in fact spurious. Modelling spatial autocorrelation can be tricky, especially for the class of binary data we will be dealing with. This ultimately required the use of software which is currently undergoing development in the R statistical package, the software of choice for this project.

The R software (R Development Core Team, 2012) contains numerous packages for fitting the binary response data we have. Our task is to experiment with different models and packages to obtain appropriate functions for fitting the kinds of models required to suit our data. As our background environmental variables are themselves highly cross-correlated, finding important associations can be difficult. Appropriate application of model selection methods is one strategy for reducing multicollinearity, although multicollinearity is less of a problem when the purpose is prediction rather than explanation (Shmueli, 2010). For this reason we will want to compare a range of model selection options.

3. Data

The data received from JNCC consisted of locations where tracked birds were recorded to be foraging; locations where birds were commuting were removed from the data prior to analysis. Foraging observations were thinned to every 10th observation prior to analysis. Environmental covariates provided are shown in Table 1; all except sediment, which is a categorical variable used to derive sand, were considered as potential covariates in the analyses.

Table 1 – Environmental covariates

Parameter	Data set	Source	Date collected	Processing	Original scale and projection	Data type
Chl_month	Chlorophyll-a concentrations, mg/m ⁻³ , monthly	PML	2009	Images taken at 1.2km square, re-mapped to 1km square	Approx 1.1km Transverse Mercator	Continuous average concentration values by month
Dist_col	Distance to nearest con-specific colony (m)	Nearest colony identified from JNCC tern colony maps	N/A	Hawths Tools distance between points in ArcGIS	1km ² grid cells OSGB 1936 Transverse Mercator	Continuous distance values (metres)
Dist_shore	Distance to nearest mainland coast (ie shortest distance to coast)	Nearest coastline identified from an Ordnance Survey high water polygon	N/A	Joins and Relates in ArcMap to store distance to closest shore for each point in the environmental layers grid.	1km ² grid cells OSGB 1936 Transverse Mercator	Continuous distance values (metres)
Sal_spring	Sea surface salinity in spring (‰)	Proudman Oceanographic Laboratory	10 year simulation	Bilinear interpolation, and Inverse distance weighted interpolation to fill in missing values near the coast	0.01 ² decimal degrees GCS WGS 1984	Continuous salinity values derived from simulation of POLCOMS
Sal_summ	Sea surface salinity in summer (‰)	Proudman Oceanographic Laboratory	10 year simulation	Bilinear interpolation, and Inverse distance weighted interpolation to fill in missing values near the coast	0.01 ² decimal degrees GCS WGS 1984	Continuous salinity values derived from simulation of POLCOMS
sst_month	Mean surface temperature by month (°C)	PML	2006-2010	Images taken at 1.2km square, re-mapped to 1km square	Approx 1.1km Mercator	Continuous average temperature values by month

Strat_temp	Surface to seabed temperature difference in summer (°C)	Proudman Oceanographic Laboratory	10 year simulation	Bilinear interpolation, and Inverse distance weighted interpolation to fill in missing values near the coast	0.01 ² decimal degrees GCS WGS 1984	Continuous temperature difference values derived from simulation of POLCOMS
Northness_1s	Seabed aspect from -1 (south) to 1 (north)	Derived from Defra digital elevation model data.	NA	Aspect function followed by transformation to radians and trigonometric cosine function, in ArcGIS Spatial Analyst	Approx. 30m ² grid cells, varies slightly across extent of data. GCS WGS 1984	Continuous values from -1 to 1.
Eastness_1s	Seabed aspect from -1 (west) to 1 (east)	Derived from Defra digital elevation model data.	NA	Aspect function followed by transformation to radians and trigonometric sine function, in ArcGIS Spatial Analyst	Approx. 30m ² grid cells, varies slightly across extent of data. GCS WGS 1984	Continuous values from -1 to 1.
Bathy_1s	Seabed depth (m below lowest astronomical tide)	Defra digital elevation model.	NA	Triangulation with linear interpolation	Approx. 30m ² grid cells, varies slightly across extent of data. GCS WGS 1984	Continuous depth values
Sediment_250	Seabed sediment/substrata	British Geological Survey (DigSBS250)	NA	Simplification of DigSBS250 Folk categories supplemented by additional data	Vector dataset GCS WGS 1984	Dominant sediment type categories: mud and sandy mud, sand and muddy sand, mixed sediments, coarse sediments, rock
Sand	Seabed sediment: 'sandy' category (1) or 'other' category (0)	Based on Sediment_250 with folk triangle.				

Slope_1s_deg	Seabed slope (° incline between adjacent grid cells)	Derived from Defra digital elevation model data.	NA	Slope function in ArcGIS Spatial Analyst	Approx. 30m ² grid cells, varies slightly across extent of data. GCS WGS 1984	Continuous slope values
Ss_currents	Shear stress: Maximum tidal force (Newtons/m ²)	Defra funded Plymouth Marine Laboratory project	June-August from 1998 to 2008.	Inverse distance weighted interpolation, derived from proWAM 12km wave model	0.003 ² decimal degrees GCS WGS 1984	Continuous values
Ss_waves	Shear stress: Maximum wave force (Newtons/m ²)	Defra funded Plymouth Marine Laboratory project	June-August from 1998 to 2008.	Inverse distance weighted interpolation, derived from POLCOMS model.	0.003 ² decimal degrees GCS WGS 1984	Continuous values
Spring_front	Probability of a frequent thermal front in spring. Ratio of strong thermal fronts to observations, averaged over all years.	Defra funded Plymouth Marine Laboratory project	June-August from 1998 to 2008.	Bilinear interpolation	Approx 1.2km ² GCS WGS 1984	Probability from 0 to 1.
Spring_frt_sd	Interannual standard deviation of probability of a frequent thermal front	Defra funded Plymouth Marine Laboratory project	June-August from 1998 to 2008.	Bilinear interpolation	Approx 1.2km ² GCS WGS 1984	
Summ_front	Probability of a frequent thermal front in summer. Ratio of strong thermal fronts to observations, averaged over all years.	Defra funded Plymouth Marine Laboratory project	June-August from 1998 to 2008.	Bilinear interpolation	Approx 1.2km ² GCS WGS 1984	Probability from 0 to 1.
Summ_front_sd	Interannual standard deviation of probability of a frequent thermal front	Defra funded Plymouth Marine Laboratory project	June-August from 1998 to 2008.	Bilinear interpolation	Approx 1.2km ² GCS WGS 1984	

4. Methodology

This section provides brief details of and justification for the methodology used during this project.

4.1 Case-Control Design

We use a case-control approach as described in Aarts *et al.* (2008). It should be noted that this is actually a use-availability design (Keating and Cherry, 2004) rather than a case-control design, since the controls represent pseudo-absences rather than true absences. Logistic regression is used to model a response variable which takes the value 1 for the observations and 0 for the control (available environment) points. The exponential function of the linear predictor is then proportional to the expected density of observations (Aarts, 2012). Warton and Shepherd (2010) demonstrate that in the case of pseudo-absences that are regularly spaced or located uniformly at random over the region, the logistic regression slope parameters (but not the intercept) converge to those of the corresponding inhomogeneous Poisson point process model as the number of pseudo-absences increases. The observed data set does not consist of all locations where terns forage, or even all locations where terns are foraging at one specific point in time; instead it is locations where terns were recorded to be foraging. The logistic regression approach models the probability that a point is a presence not a pseudo-absence. This probability has no physical meaning and tends to zero as the number of control points increases; it is the intensity of the presences rather than the probability of occupancy that is of interest and a Poisson point process model therefore has a more natural interpretation. Existing software does not allow a spatial Poisson point process modelling approach to the current problem to be fitted routinely. We therefore take the approach of generating control samples and using logistic regression to approximate the point process model. However, future developments in the INLA (integrated nested Laplace approximation) package (INLA, 2012) should allow the point process model to be fitted directly without the need to generate control samples.

As we wanted to consider a range of modelling options, we proposed generating the control samples so that each set resembled the foraging locations in the original case data (particularly in relation to exhibiting autocorrelation). The simplest way of implementing this seemed to be to generate sets of initial “track starts” – one for each track in the case data – and then applying the rest of the case track movements between foraging locations to each. If a control “track start” fell on land, it was replaced with another random starting location; subsequent points falling on land were simply omitted. This, in effect, creates a set of control samples which are the original case tracks relocated randomly throughout the suggested range of the species – see Figure 1. Having control tracks (rather than just control points as per Aarts *et al.*, 2008) allowed us to think sensibly about using different forms of random effect structure when modelling. However, given the approach eventually chosen did not include random effects, future work could revert back to using control points rather than control tracks; these will provide better spatial coverage per sample point than tracks (therefore reducing the sample size and processing time needed), but will not enable small-scale assessment of spatial correlation.

Control locations (the initial track starts) were generated by taking a random angle (uniformly distributed between 0 and 2π radians) and then a random distance from the colony (also uniformly distributed between zero and the maximum foraging range). This gives a greater density of control points closer to the colony than if points were uniformly distributed over a circle. Aarts *et al.* (2008) recommend that control points should be selected in proportion to accessibility; this would certainly mean a greater density closer to the colony, although it does not necessarily imply a uniform distribution of distance. Selecting points in proportion to accessibility means that the model outputs

provide estimates of preference, which is defined as the ratio of use to availability. Estimates of usage can then be obtained by multiplying preference by accessibility. On the other hand, if the control points were uniformly distributed over the circle this would mean that the model outputs would provide direct estimates of usage. The number of control points needs to be sufficiently large to ensure that estimates of slope parameters converge to stable estimates. Exploratory analysis by JNCC suggested that up to five or six control tracks may be required. To be cautious, twelve replicate control tracks were generated for each case track. Further research is needed to determine the number of control points required to achieve stability for different sample sizes as this will not necessarily be a multiple of the number of observations. This is likely to be a complex problem, as the number of control points required could potentially vary hugely between colonies – it will be a function of colony size, colony density, topography of nearby land masses, interactions with other species, availability and distribution of food sources, and so on.

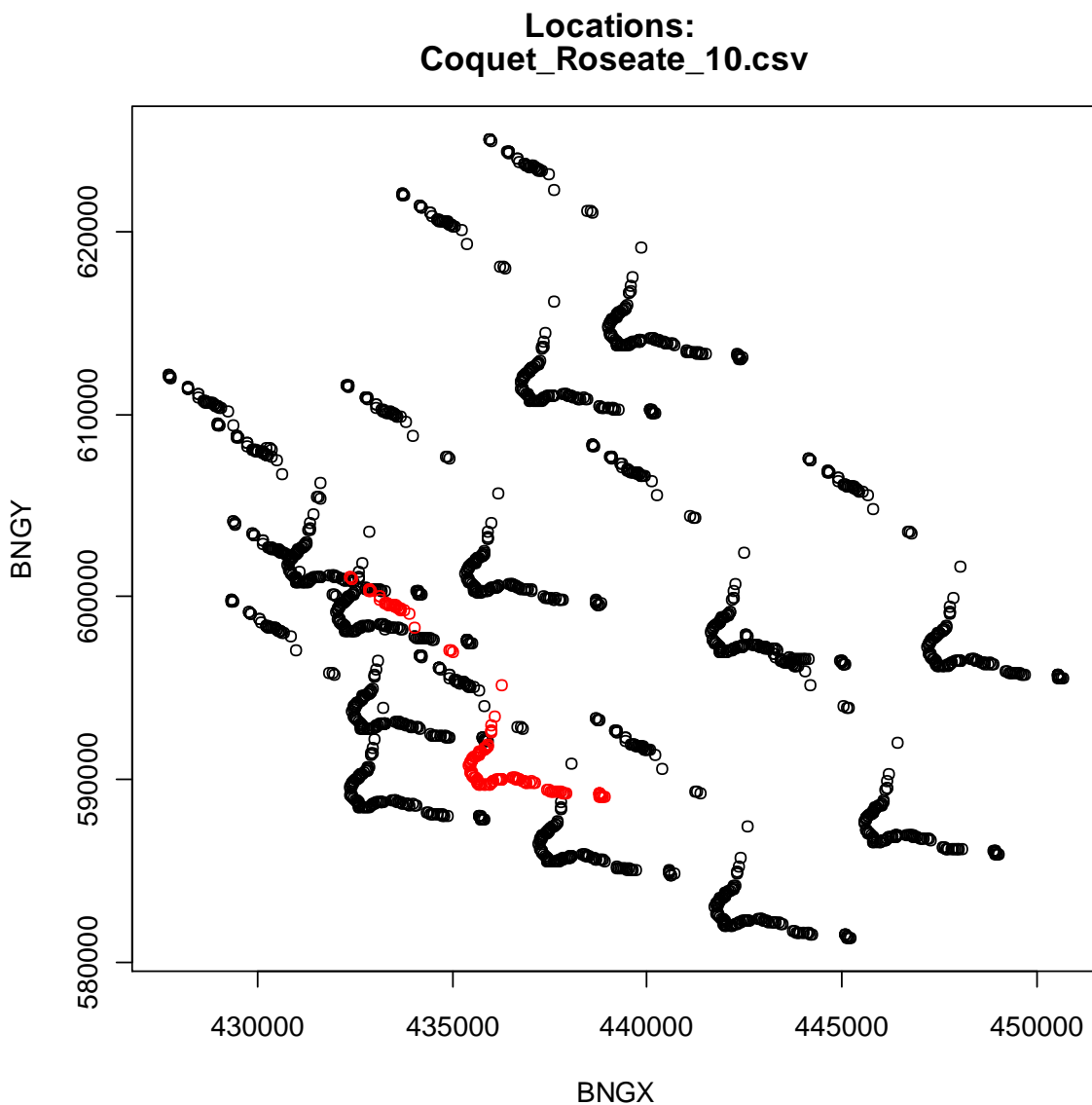


Figure 1 – Example (for roseate tern at Coquet) of the pseudo-absences (in black) generated by shifting the original track (in red)

The environmental information described in Table 1 was attached to each case and control point. To make grids for interpolating the preferences (and usage) to the whole foraging range of each species at each colony, environmental information for points spaced every 500m were provided for the area within the maximum foraging range for each species, around each colony. This was either 30km (Arctic, common and roseate terns) or 60km (Sandwich terns).

4.2 Statistical Model and R Packages

The modelling needs to take account of the fact that points on a given track are not independent (i.e. they are repeated measures), as birds generally move only a short distance from one observation to the next. Failure to account for the lack of independence between observations within track – a form of pseudoreplication – leads to underestimation of the variance of parameter estimates and might therefore result in some environmental variables being wrongly regarded as significant. The simplest way to deal with this is to weight each observation by the reciprocal of the length of the observed track. This has the effect of treating each track as a single sampling unit, instead of treating the individual observation as the sampling unit, but is preferable to averaging data over tracks as individual covariate values are retained for each point. Thus, the basic statistical model used in our analysis is a weighted binomial generalised linear model (GLM) or generalised additive model (GAM) with a logit link. As the data are binary the dispersion parameter should be fixed at 1.

Existing methodological approaches (for example, Wakefield *et al.*, 2011) fit a generalised linear mixed model (GLMM) to account for between-individual variation, although it is unclear exactly how these authors deal with the control points when specifying the random effects. We considered several approaches to the tern analysis using random effects. The standard approach to repeated measures would be to fit a separate random effect for each case or control track. However, this is not appropriate in our situation, because the case tracks consist entirely of presences while the control tracks consist entirely of absences, leading to biased estimates of the fixed effects of the covariates. An alternative approach that we considered, which was probably the one used by Wakefield *et al.*, is to fit a random effect for each combination of a case track with its corresponding control tracks. Unfortunately, this approach is not suitable either, as within each level of the random effect there are units from the same control set that are highly correlated and units from different control sets which are uncorrelated; the model incorrectly assumes that these correlations are equal. For smaller numbers of control sets many of the entries in the correlation matrix are between units from the same track, and as in the previous approach parameter estimates are biased. It is likely that a variation of such a model accounting for spatial autocorrelation would help reduce this bias. As the number of control tracks increases, so entries in the correlation matrix that correspond to units in different control sets start to dominate; this means that the variance of the random effect tends towards zero as the number of control sets increases and the results approach those for an unweighted GLM that does not account for dependence within a track/control set combination. As described above this gives unbiased estimates of the parameters, but estimates of the standard errors are biased (i.e. the parameter estimates falsely precise) because they fail to take account of the dependence within a track. Under-estimation of the standard errors would have been even greater if we had used control points instead of control tracks.

We felt that none of the approaches we explored for specifying the random effects were appropriate, thus we chose to use a fixed effects only weighted model. Weighting avoids biasing the results towards the longer tracks, but would still be necessary even if all the tracks were of the same length because of the autocorrelation between observations in the same track. As there is not time for birds to fly far between successive foraging observations and as the environmental covariates are spatially autocorrelated, observations from the same track tend to occur in similar environmental conditions. This does not necessarily indicate differences in individual preference between birds; multiple tracks per bird would be needed to investigate this.

Our model exploration work led us to three different functions within R, representing different levels of complexity, but all fitting variants of the same basic model:

1. Function `glm()` in the base “stats” package;
2. Function `gam()` in the “mgcv” package (Wood, 2011); and
3. Function `inla()` in the developmental INLA package (INLA, 2012).

All three functions can fit the weighted binary logistic regression required. The first function, `glm()`, is the GLM workhorse of R, is extremely stable and reliable and has a large suite of generic helper functions. The function `gam()` in package `mgcv` provides a simple, flexible interface for fitting smooth functions for the effect covariates, in case some relationships are more complex than straight lines. Finally, we looked for a suitable method to account for residual spatial autocorrelation (i.e. the remaining autocorrelation after the effect of the environmental covariates has been accounted for). Not accounting for this autocorrelation can lead to underestimation of the standard errors for parameter estimates, and this in turn implies some variables may wrongly be declared significant as a consequence; for more detail see Beale *et al.* (2010). The track structure of the data means that a simple rectangular grid would not be appropriate, and the data sets are too large, in general, for modelling large covariance matrices based upon distances between observations (note standard methods consider all pairwise differences). The `inla()` function in the INLA package solves this problem by defining a “mesh” to represent the spatial autocorrelation structure based upon the locations of the observations; in this way, one observation is linked to a small number of its closest neighbours, and the modelling is intelligent enough to account for these between-neighbour distances. Happily, this method has allowed us to fit full spatial models with whole-species data sets. For an example of the mesh structure produced by the INLA package, see Figure 2.

In summary, the analysis takes three distinct stages; firstly, we fit generalised linear models and perform model selection. The high level of correlation between explanatory variables suggests that we need to be careful about model selection and use different methods depending on whether the aim is prediction or explanation. Because of this concern, we in fact use three different approaches to selection on the basis of: AIC (Akaike’s Information Criterion); BIC (Bayesian Information Criterion); and significance of individual terms via likelihood ratio tests (LRTs). The AIC is known to produce large models, with many explanatory variables which may not be statistically significant, but is optimal in terms of linear models for prediction. The BIC contains a penalty (against adding variables to a model) which is stronger than that for AIC and which is a function of number of observations. BIC is consistent in the sense that if the true model is among the candidate models then the probability of selecting it approaches 1 as the sample size increases, but it is sub-optimal for prediction. As our “number of observations” (which includes the pseudo-absence data) is typically very large, the penalty is itself very large and hence the models selected by BIC tend to be much smaller than those chosen by AIC.

Constrained Delaunay triangulation

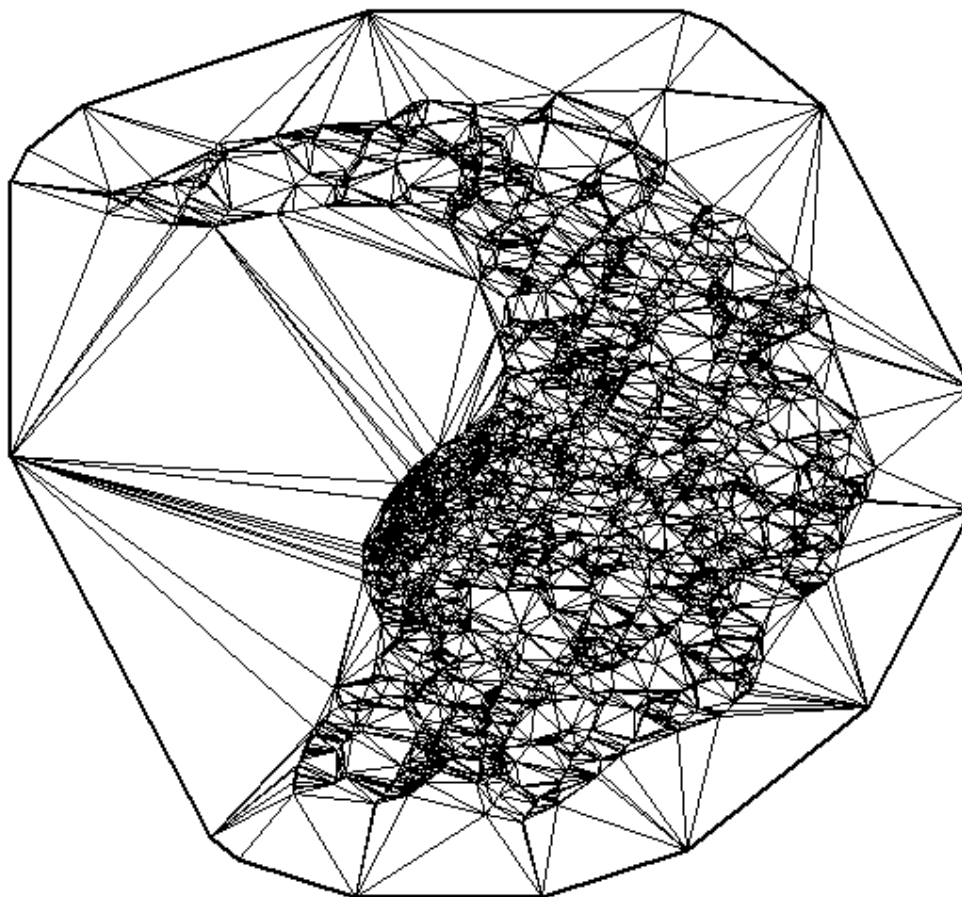


Figure 2 – Mesh to represent spatial correlation structure of response data; example shown for Forvie colony and Sandwich terns

The LRT method is preferred for identification of important associations between tern presence and the covariates, as each term in a selected model must be individually statistically significant; we use a significance level of 10% here as we wish to explore feasible models thoroughly. Note, however, that if two covariates are highly correlated and both are included in the model, neither may appear significant, even though either one of them would be significant if fitted individually. Even in the absence of correlation between the potential covariates association does not necessarily indicate causality as the significant variables may be correlated with a missing, causal, covariate. At this point we also check for consistency of covariate effects across different years, for those species/colony combinations which have more than one year's worth of data; in essence, we fit interaction terms between year and the LRT-selected covariates and conclude effects are not temporally consistent if the interaction is significant.

These generalised linear models will not account for non-linear relationships between response and covariates. For this, we fit generalised additive models (GAMs) which use “spline functions” to

describe these response curves – these are non-linear, non-parametric functions which can take any shape, limited only by a number of “basis functions” which control how “wiggly” the line is allowed to be. We chose three basis functions so that the response-covariate relationships would be fairly smooth – we do not believe *a priori* that the relationships should be too complex. For GAM models, as the AIC is only approximate, model selection was instead performed by assessing the significance of smooth terms by comparison to a chi-squared distribution. Next, the INLA method is used to assess the effect of any spatial autocorrelation in the response data, over and above that explained by the environmental covariates in the model. INLA is applied to the model chosen by the GAM model selection; most models chosen across the species/colonies contained some smooth terms, and we wished to have a consistent approach. INLA allows us to account for spatial autocorrelation within the model and to identify those parameters which may have been falsely identified as significant.

Finally, we use the fitted models to produce maps of the estimated preferences and estimated usage, as detailed in Aarts *et al.* (2008). Again, for reasons of consistency we aimed to use the selected GAM model where possible, although variables which the INLA model suggests are no longer significant after allowing for spatial autocorrelation could be removed. The maps of preference show which “habitats” the terns might prefer, to a degree, but note this interpretation is not strictly applicable if either “distance to colony” or “distance to shore” are in the selected model.

If there was no preference for particular habitats, the odds ratio would be equal to the ratio of the number of observations to the number of controls. To calculate preference the odds ratio is adjusted by multiplying by the number of controls per observation, which is equivalent to taking the exponential transformation of the linear predictor and multiplying by the number of controls per observation. This does not purely reflect preferences for environmental variables because of bias in our accessibility model, which assumed a uniform distribution of distance from colony. To correct for this the linear predictor could be broken up into two parts: a component due to distance from colony and a component due to environmental variables alone.

The maps of usage adjust for accessibility, using the same model as that for generating the control sample locations – i.e., a uniform distance distribution from the colony. Equation (5) of Aarts *et al.* is:

$$f_u(X_s) = \frac{h(X_s)}{1 - h(X_s)} r f_a(X_s)$$

where $f_u(X_s)$ is the spatial probability density function for usage, $h(X_s)$ is the predicted value from the fitted model, r is the number of controls per observation and $f_a(X_s)$ is the probability density function for accessibility. In our case

$$f_a(X_s) \propto 1/\text{distance}$$

so preference is divided by distance to colony and multiplied by a scale factor which ensures that the probabilities sum to one. These results could then be multiplied by the number of birds in the colony.

4.3 Verification of the Weighted Model

How do we know the weighted regression model is appropriate? This is easily verified by thinking about what the model would look like if there were only a single observation per bird, rather than the tracks that we have currently. There would then be no basis for random effects or weights, and we could fit a simple logistic regression (preferably accounting for spatial autocorrelation in some way – but that is not relevant to the point we’re trying to make so we’ll ignore it for now). We have tested the “single point” model by sampling randomly a single point from each track, and then running this simple analysis. By generating multiple single sets, we were able to gauge how variable the “single point” results were. Although the results will be quite variable, they provide the perfect yardstick for the “full track” analysis. The weighted logistic model parameter estimates were very similar to the mean estimates from the single point analysis and the standard errors were also similar, as was to be expected, since the weights are effectively accounting for the track lengths; in addition the results were very insensitive to the number of control sets used.

Note that the “single point” analysis represents the maximum amount of thinning that could be applied to the data. In consequence, the single-point analysis described above also (at least partially) resolves the issue of thinning; as long as the data are not thinned too much – where “too much” would imply that possible relationships with covariates are lost, perhaps by obtaining a track on a coarser scale than the covariates – the actual level of thinning is unimportant.

4.4 Other Modelling Approaches / Packages Explored

When the use of a mixed effects modelling approach was investigated initially (before it was found unsuitable), a variety of functions /packages were explored. For fitting a GLMM, both `glmer()` in package `lme4` and `inla()` again appeared good possibilities. The function `glmmadmb` in package `glmmADMB` did not perform well and, in any case, does not appear to allow a weights option. Function `glmmPQL` in package `MASS` was not suitable for two reasons: the PQL (penalised quasi-likelihood) fitting process is known to produce biased results; and there does not seem to be a possibility of fixing the dispersion parameter. The `gamm()` function in `mgcv` appears to use PQL and therefore suffers from the same problems as `glmmPQL`.

4.5 Discussion

Our tender suggested exploring a mixed models approach, as this appeared to have been successful in earlier work (e.g. Aarts *et al.*, 2008; Wakefield *et al.*, 2011). In Aarts *et al.* (2008) there are multiple observed tracking instances of the same bird – this would seem to be the difference between their data and the current terns data, and may explain why a mixed model approach may have worked for them and not for us. However, it also appears that these earlier studies may have ignored the fact that whereas observations within a track are serially correlated, there is no correlation between an actual observation and the corresponding control observations. This may have led to standard errors being under-estimated in these studies. We might speculate that a weighted analysis may have been beneficial in these earlier works as well as with our current analysis.

5. Results

This results section contains a summary discussion of the results obtained for the analyses of the twelve data sets. This section is ordered by colony within tern species. The full R output for each model is contained in the “Results Appendix” accompanying this report. The sections in that document mirror those here – for example, section 5.1.1 below summarising the Coquet Arctic tern results is mirrored in section A.1.1 in the Results Appendix document. This should enable straightforward cross-referencing where required.

Note that the models presented are the ones chosen via model selection procedures. Owing to the very large correlations between the environmental covariates, it would be wrong to infer *causal* relationships from the *correlative* relationships discovered; we are finding association rather than proving what is driving habitat selection. Nonetheless, this does not preclude obtaining good predictions. In other words, association rather than causation can provide the information we require to assess which areas are suitable for foraging behaviour.

5.1 Arctic Terns

5.1.1 Coquet

This analysis illustrates the harsher penalty from the BIC for model selection; that produces a model containing only distance to colony, whereas AIC selects distance to colony, sst_may, summ_front_sd, ss_wave and ss_current. LRT selects the same model as AIC, and the effects across the three years’ worth of data were consistent.

The GAM model selected also uncovered distance to colony, sst_may, summ_front_sd, ss_wave and ss_current; an illustration of the relationships is shown in Figure A.1.1.4. From this it can be seen that preference increases with increasing values of sst_may and ss_current, and decreases with increasing values of ss_wave and summ_front_sd, although the relationship with sst_may seems to be largely driven by a small number of outliers. Of those four, the relationship with ss_wave seems strongest ($p=0.0125$). The relationship with distance to colony shows that preference decreases steadily as you move further away from the colony (to about 20km) and then the relationship levels off. This is the strongest effect of all ($p<0.001$).

The INLA model suggests that spatial autocorrelation may be causing the suggestion of significance for some of the variables noted above; from the INLA output, we have credible intervals rather than p-values, and we regard a variable as important if the interval does not overlap zero. In fact, only sst_may (0.0994, 3.2528), ss_wave (-0.0804, -0.0129) and the smooth distance to colony relationship remain as important covariates once spatial autocorrelation is taken into account. However, the smooth relationship for distance to colony contains a seemingly unlikely upturn for large distances, likely an artefact caused by a lack of case data at such distances (see Figure A.1.1.5).

5.1.2 Farnes

Here AIC selects distance to colony, sst_april and summ_front_sd, while BIC leaves out sst_april of the three. LRT again agrees with AIC. With only one year’s worth of data, there was no need to check for consistency across years. GAM also selected distance to colony, sst_april and summ_front_sd, and as seen in Figure A.1.2.4 the relationships with distance to colony and sst_april

are decreasing whereas that with `summ_front_sd` suggests higher values correspond to higher preference. For this data set, the INLA analysis supports the model selected as none of the credible intervals overlap zero. The covariate relationships were sufficiently linear for INLA to dispense with the GAM element.

5.1.3 Outer Ards

AIC selects distance to colony, distance to shore, `chl_apr`, `chl_jun` and `sst_may`; BIC on the other hand finds only `chl_apr` and `sst_may`, ignoring the distance variables completely. The LRT selection adds distance to colony to the BIC selection. There was no evidence of covariate effects being different across years. Interestingly, the GAM model selects distance to colony, `chl_apr` and `sst_april`, differing from the LRT selection by switching `sst_may` and `sst_april`. Figure A.1.3.4 shows that for all three selected GAM covariates, higher values are associated with lower preference. The spatial INLA analysis provides three credible intervals, none of which overlap zero, so there is no evidence on this occasion that any of the covariate effects found are due solely to spatial correlation. Again, the covariate relationships were linear and so INLA ran a spatial GLM rather than a spatial GAM.

5.2 Common Terns

5.2.1 Coquet

Here AIC selects distance to colony, `chl_june`, `sst_april` and `ss_wave`, but as with the Arctic terns, BIC selects only distance to colony. The LRT model selection matches the AIC choice, and note here there is very slight evidence of the effect of distance to colony varying across years; however, we note that $p=0.07443$ and hence the evidence is weak; we assume the effects are consistent across years, therefore. The GAM model chosen contains the same covariates as the AIC and LRT selections. Figure A.2.1.4 suggest two linear relationships, with `chl_june` (increasing) and `ss_wave` (decreasing), and two non-linear relationships. With distance to colony, the relationship is clearly decreasing, but the rate of decrease increases further away from the colony. The relationship with `sst_april` is also decreasing, but perhaps levelling out for the upper end of the range of temperatures. The spatial INLA analysis does not contain any intervals overlapping zero, so this reaffirms the significance of the four selected variables. The plot for `sst_april` in Figure A.2.1.5 is interesting; it suggests that the negative slope is due to the outlying values of temperature at the bottom end of the range, and that in fact the relationship is an increasing one for the bulk of the data in the region around 7 °C. Further work should consider re-running the model excluding these outliers.

5.2.2 Larne Lough

AIC selects distance to colony, distance to shore, `sst_april` and `bathy_1sec`. BIC finds three of these, dropping `sst_april`. LRT and GAM agree with AIC. There is no evidence of these covariate effects being different across years. The relationship with distance to colony is negative, whereas for distance to shore and `sst_april` the relationships are positive (see Figure A.2.2.4). The relationship with `bathy_1sec` is “cup-shape”, having lowest preference in the -150 to -100 metres range, but higher when deeper than 150 metres or shallower than 100 metres. The spatial INLA analysis supports the importance of all the variables except `bathy_1sec`; note the very wide interval estimates in Figure A.2.2.5.

5.2.3 Leith

For Leith AIC selects a large number of variables: distance to colony, distance to shore, chl_may, chl_june, sst_may, summ_front, spring_front, sal_spring, bathy_1sec, slope_1s_deg. However, the large size of the data set means that the penalty on adding variables with BIC is huge, and therefore only distance to colony is selected. The selection with LRT is large, but not quite as large as for AIC; the variables here are distance to colony, distance to shore, chl_may, chl_june, sst_may, sal_spring, bathy_1sec and slope_1s_deg. With so many variables, it is not surprising that one (sst_may) displays weak evidence ($p=0.0604$) of varying between years, but again we choose to assume a consistent effect across years. (Note that the consequence of not doing so means that we can then only make predictions on a per-year basis, and only for the years for which we have data.)

The GAM model selected matches the LRT model. It can be seen clearly in Figure A.2.3.4 that most of the variables selected have quite small effects on preference relative to the effect of distance to colony, by far the strongest effect ($p<0.001$). The variables chl_may, chl_june, sst_may and sal_spring all have very small positive relationships with preference. The effect of slope_1s_deg is a little greater, although there is uncertainty for high values owing to the relative sparseness of data. Distance to shore shows a moderately negative association with preference. There is a slight non-linear relationship for bathy_1sec, with an increasing relationship for lower depths levelling out for higher values. Distance to colony has a negative relationship which tends to get steeper for higher distances.

The spatial INLA analysis suggests the relationships with chl_may and sal_spring may simply be due to spatial correlation in the response data, as their credible intervals overlap zero, $(-0.0007, 0.0314)$ and $(-1.6119, 2.0057)$ respectively.

5.2.4 Mull

For Mull, the AIC selected chl_apr, chl_may, sst_april, sst_may and ss_wave. BIC however selected only chl_apr. LRT was different again, finding chl_apr, chl_may and ss_wave. And GAM was different too – the variables here were chl_apr, chl_may and bathy_1sec. It should be noted here that the Mull colony appears to be surrounded by a lot of land on all sides, with the sea (or lochs) forming narrow strips around the colony. It is not too surprising in particular that distance to colony is not significant – there are areas of water close to the colony which have land in between; it may not have always been possible to continue tracking terns which flew over such land, hence a number of possible areas close to the colony “as the tern flies” may not have been sampled, and thus may have caused bias in the results. We do note that chl_apr and chl_may do occur consistently in chosen models.

Figure A.2.4.4 suggests that all three relationships from the GAM model are non-linear. For both chl_apr and chl_may, preference increases with increasing chlorophyll for low concentrations, then levels off for higher values. For bathy_1sec the relationship is generally negative, levelling out or even increasing slightly for higher values (shallower depths). The spatial INLA analysis broadly supports the GAM model.

5.3 Roseate Terns

5.3.1 Coquet

The AIC chooses a large model here; the variables are distance to colony, distance to shore, chl_may, sst_june, sst_may, summ_front_sd, strat_temp and ss_current. The BIC model is much smaller, having only distance to colony, sst_june and sst_may. The LRT model is in between, being composed of distance to colony, chl_may, chl_apr, sst_june, sst_may and summ_front_sd, and there was no evidence of differential year effects. The GAM model selected matches the LRT model. As shown in Figure A.3.1.4, the effects are all linear, with the strongest being for sst_june and sst_may (both $p < 0.001$), although the relationship is positive for sst_may but negative for sst_june. There are weaker positive relationships with chl_apr and summ_front_sd, and weaker negative relationships with distance to colony and chl_may.

For the spatial INLA analysis, the credible interval for chl_apr is $(-0.0034, 0.6158)$, which (only just) overlaps zero, casting doubt on the importance of that variable in terms of association with preference.

5.4 Sandwich Terns

5.4.1 Coquet

For Coquet, the AIC chose distance to colony, distance to shore, chl_june and sst_april. BIC only chose the two distance variables. LRT also chose distance to colony, distance to shore, chl_june and sst_april, and spotted a potential cause for concern re: interactions with year – the interaction of distance to shore and year was marginally significant with $p = 0.04994$. Again we choose to ignore this interaction, so the maps produced (and presented in Section 7) work on the basis of average effects across all years (produced from the model without interaction term).

The GAM model matches the AIC and LRT models. Figure A.4.1.4 shows the four relationships, which all have different characters. The relationship with chl_june is very small and positive. That between preference and sst_april is non-linear, being flat for lower temperatures but decreasing for higher temperatures. There is a reasonably strong negative relationship with distance to shore, but the strongest effect of all (if not the most significant) is for distance to colony, generally decreasing, but faster with longer distances. Note that it is the high uncertainty with the distance to colony effect – characterised by the wide interval estimates in Figure A.4.1.4 – that cause distance to colony to have a higher p-value than distance to shore, despite having ostensibly a stronger relationship with preference.

With the spatial INLA analysis, the credible interval for chl_june overlaps zero – $(-0.0439, 0.1972)$ – which suggests it is not really important, and was only found to be in the GLM/GAM analyses due to spatial autocorrelation in the response data.

5.4.2 Farnes

Here the AIC selects distance to colony, distance to shore, summ_front, spring_front and ss_wave. BIC picks only distance to shore. LRT agrees with AIC, but the GAM model chosen drops distance to colony, despite the four remaining variables displaying linear relationships. As can be seen in Figure A.4.2.4, summ_front and ss_wave have small positive relationships with preference, whereas

spring_front has a small negative one. The strongest relationship (and with $p < 0.001$) is for distance to shore, very strongly negative, suggesting the birds like staying closer to shore.

There is no evidence in the credible intervals from the spatial INLA analysis that spatial autocorrelation influenced the model selection.

5.4.3 Forvie

For Forvie, the AIC selects distance to shore, sst_june and strat_temp. BIC drops sst_june, but the LRT selection matches that of AIC. The GAM selection agrees with BIC, composed of distance to shore and strat_temp. The effect of both of these variables is strongly negative (see Figure A.4.3.4). The spatial INLA analysis does not cast doubt on the earlier model selection.

5.4.4 Larne Lough

This data set proved very problematic; there were so few combinations of covariates for the case data, it was only feasible to use a very limited set of covariates. In fact, we chose distance to colony and distance to shore, as these were most likely to be important and were at least likely to be measured accurately. Both variables were indeed selected by *all* the different criteria. There was no evidence of the effects varying between years, and the spatial INLA analysis did not provide credible intervals overlapping zero.

From Figure A.4.4.4, it can be seen that the relationship between preference and both distance to colony and distance to shore is negative, and that the effect of distance to shore is stronger (with $p = 0.002$).

5.5 Overall Comments

Looking across the colonies for each species, we find a number of consistent features. Note, however, that these are only associations; we cannot state that the selected covariates are causal drivers of foraging behaviour. Of course, it may be that the true drivers are in fact unrecorded.

For Arctic terns, distance to colony is an important predictor, and each analysis unearths one of the sst_ variables. Note that owing to correlations between the covariates, it is not possible to disentangle these effects further. Distance to colony is of course related to accessibility, but it is noticeable that most analyses find significant effects over and above this variable.

The situation for common terns would seem to be similar to that for Arctic terns, in that both distance to colony and an sst_ variable are selected consistently. Mull is an exception here, but as noted earlier the geography surrounding that colony is problematic, in that it has likely caused biases in the data collected.

There is only one colony for roseate terns, which adds chl_ variables to the distance to colony and sst_selections from the previous two species.

For Sandwich terns, the only consistent explanatory variable chosen was distance to shore; the individual colonies mostly produced other significant effects, but these varied quite widely.

6. Discussion

The methodology used presents an advance on earlier methods, with the weighted regression accounting for the differing numbers of observations in each of the recorded tracks. This method is different from that developed by Aarts *et al.* (2008), which does not explicitly use weights and hence implicitly ascribes more weight to tracks with more locations along them. We have made this alteration because the high correlation between locations on the same track makes track a better sampling unit than location

The maps in Section 7 illustrate clearly that the models describe some but not all of the predictability for preference of search/foraging area. With some maps – for example, the common tern for Coquet, the Sandwich tern for Forvie – the red shaded areas (representing higher preference) coincide well with the observed data. For others, however, the observed data lie on a mixture of red and blue areas, suggesting that the given explanatory variables are not sufficient to explain well the preferred foraging areas for those species/colonies; examples here include the common tern for Larne Lough or the Arctic tern for Outer Ards. However, note that we would expect to see some observations in the low preference areas as they are spatially more extensive.

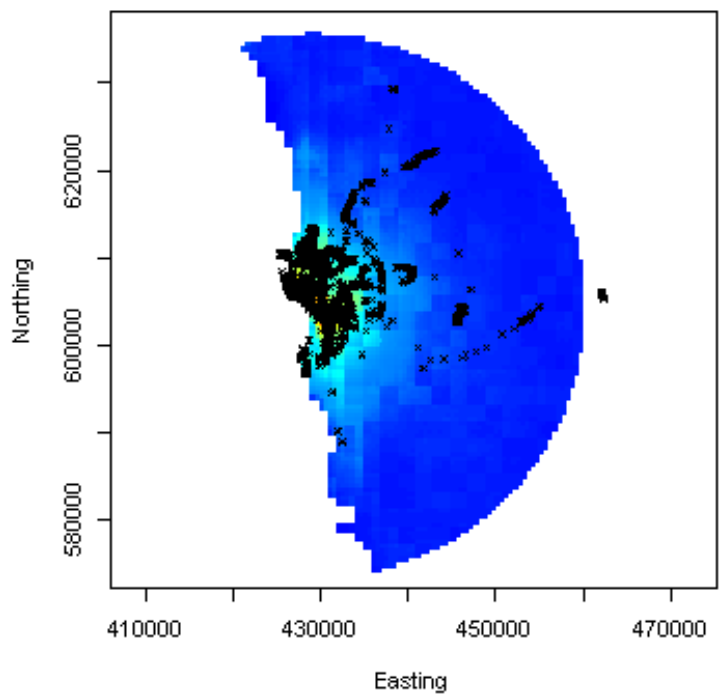
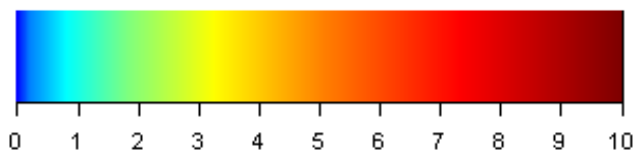
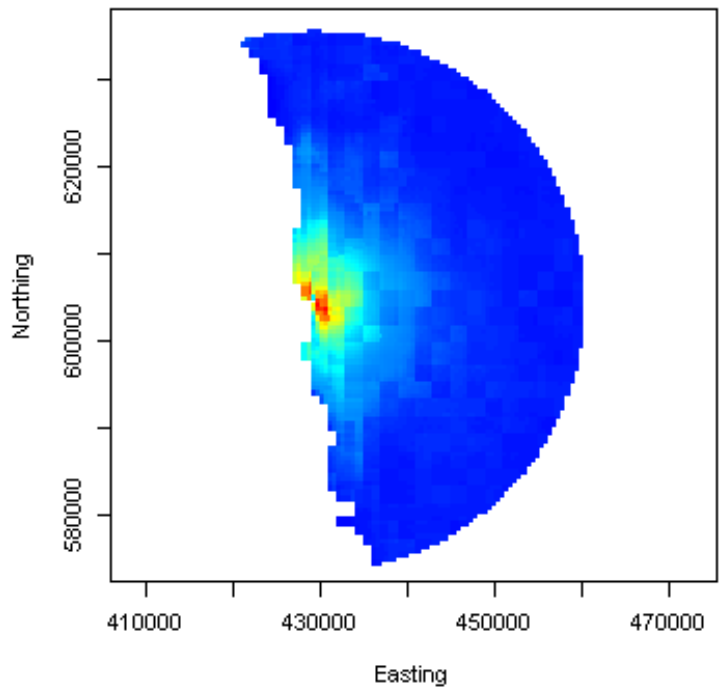
Some of the large sample sizes effectively mean the analysis should be very sensitive to genuine relationships between tern presence and the environmental covariates. However, there is also a real possibility of uncovering spurious relationships, that is, relationships which may be due to a very small actual effect, caused by only a small number of observations that are outlying in covariate space. Further work should investigate the effect of removing outliers in the environmental covariates that may be unreliable. Despite this caveat, there was some evidence of consistent effects between colonies within species of tern. It may be possible that such differences within species but between colonies are due to the birds reacting to local variability. Future work could look at the extent to which it is possible to make reliable predictions for one colony using a model developed for another colony.

7. Maps Appendix

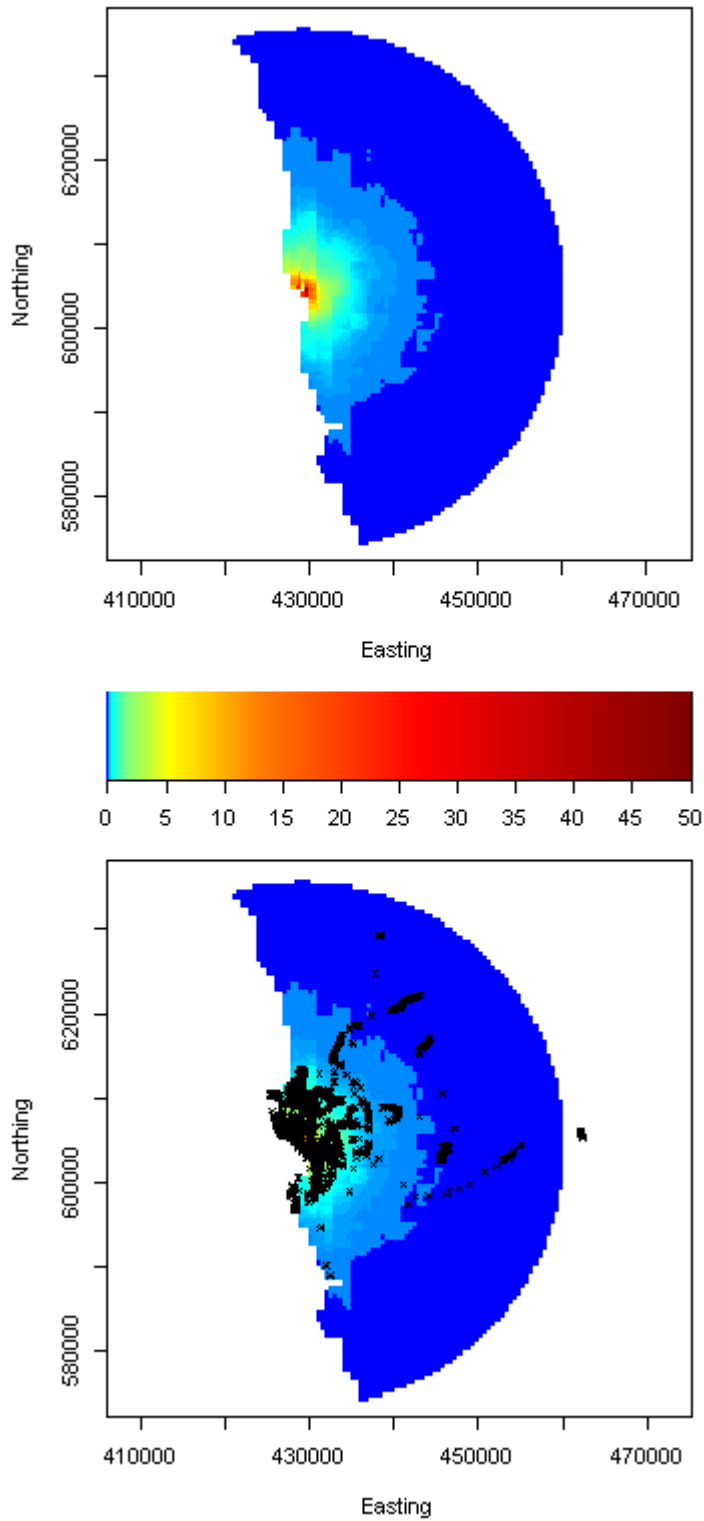
The maps in this appendix have been produced on the basis of the selected GAM models described in Section 5 and in the Results Appendix. Some outliers in the covariate values were removed before making predictions. For each species/colony combination, there are two pages of maps. The first page contains mapped preference, the second mapped usage. For ease of comparison with actual observations, each predicted map has been plotted twice – once by itself (top) and once with the observations added as black crosses.

For the usage maps, the probabilities have been multiplied by 1000.

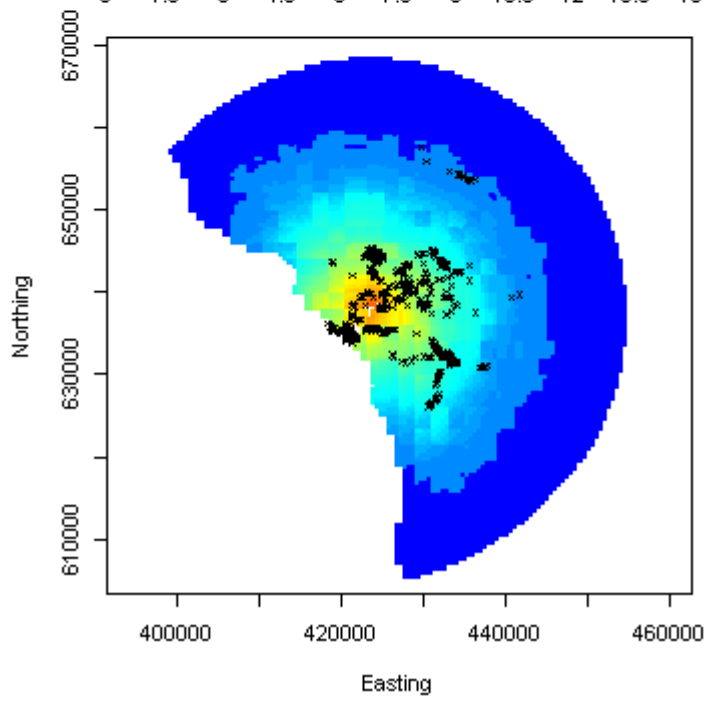
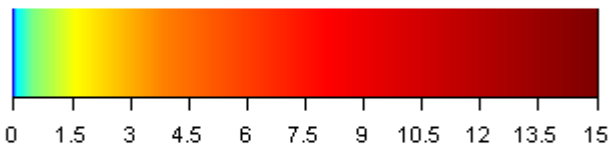
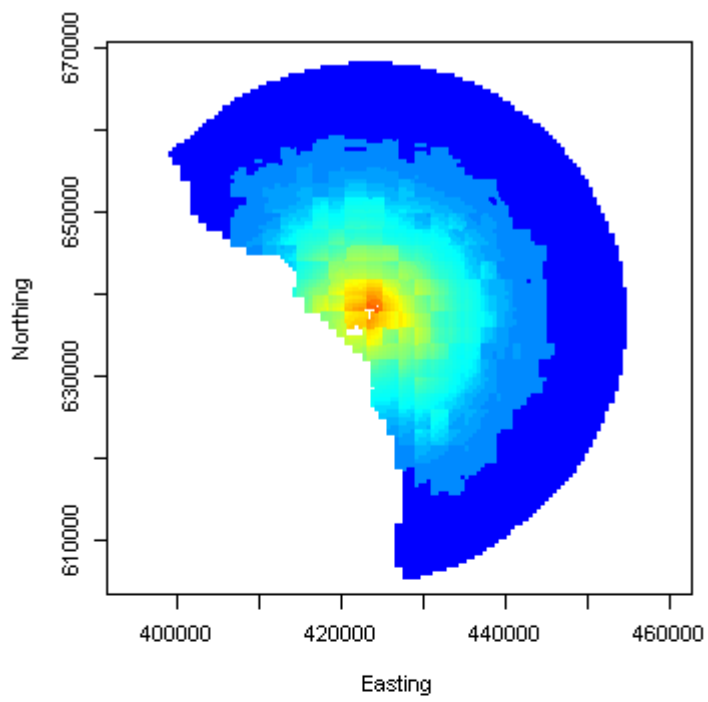
Arctic Tern, Coquet Colony – Preference



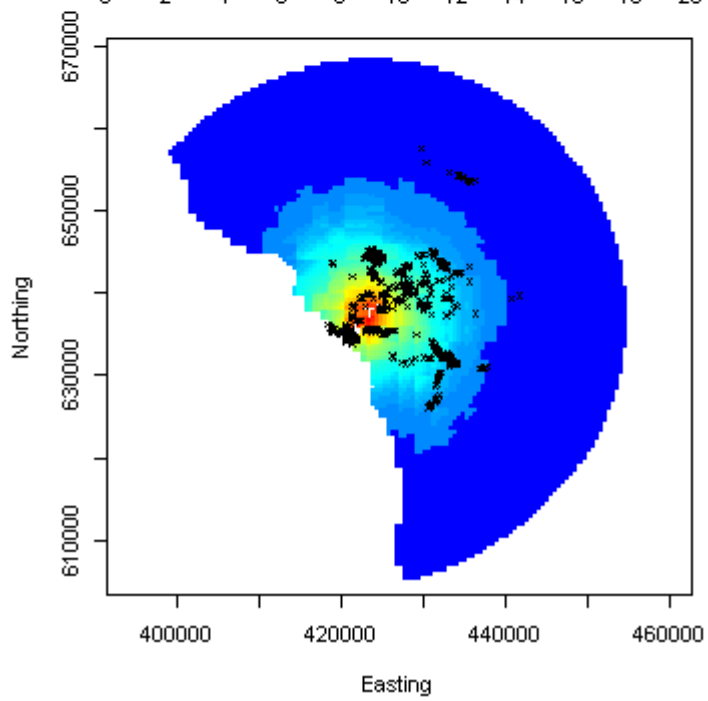
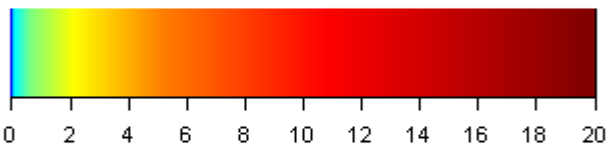
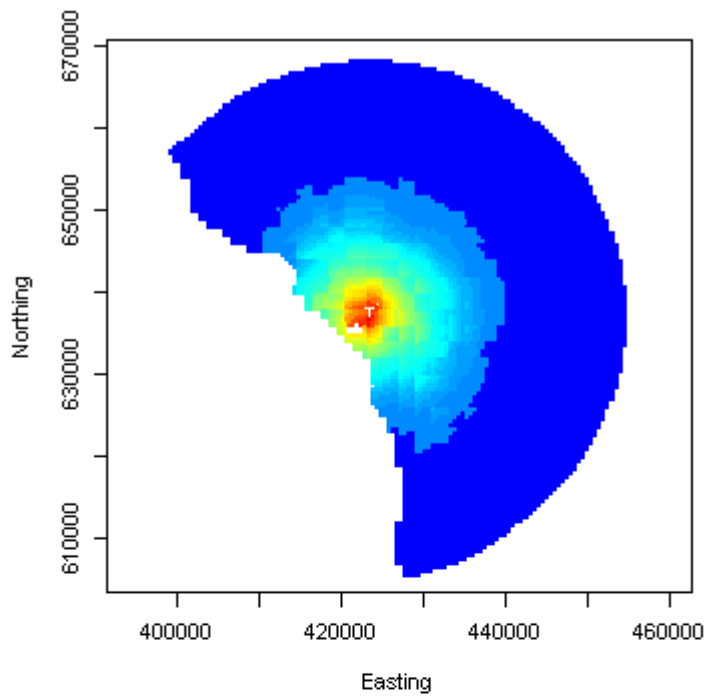
Arctic Tern, Coquet Colony – Usage



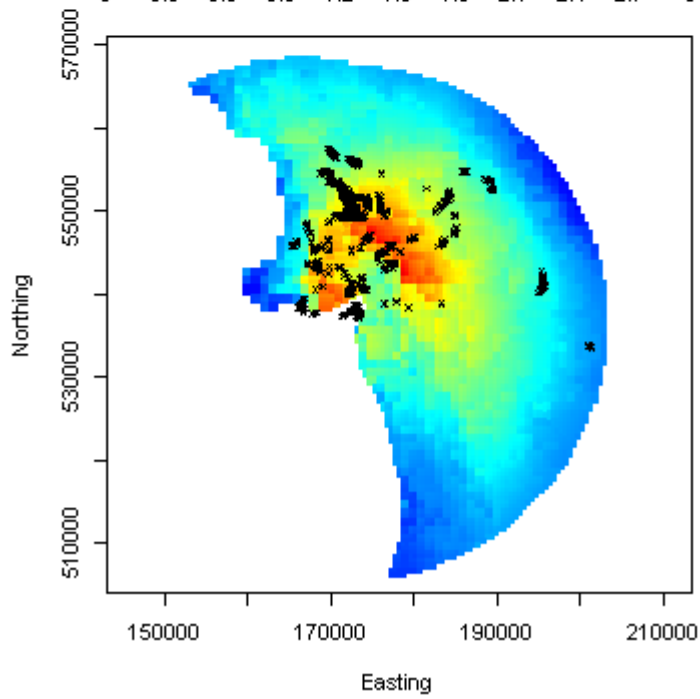
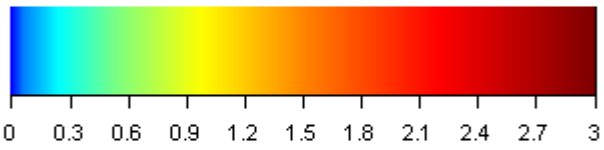
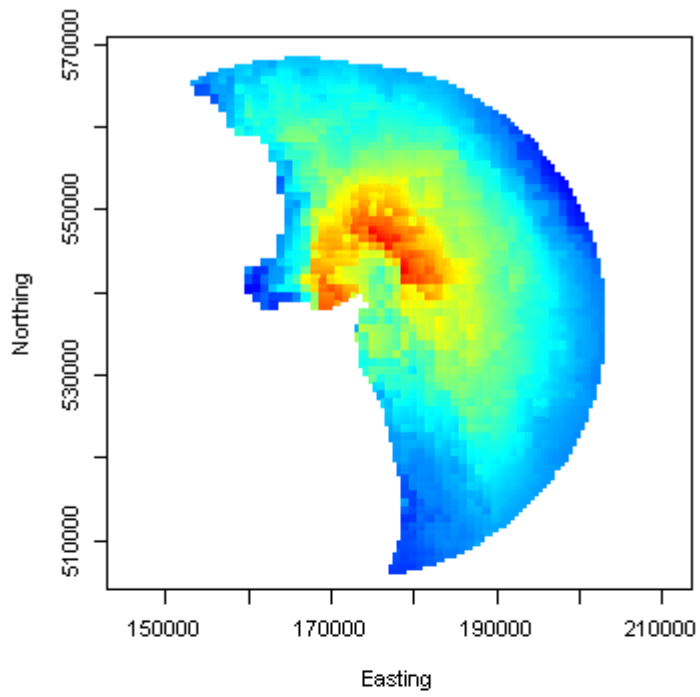
Arctic Tern, Farnes Colony – Preference



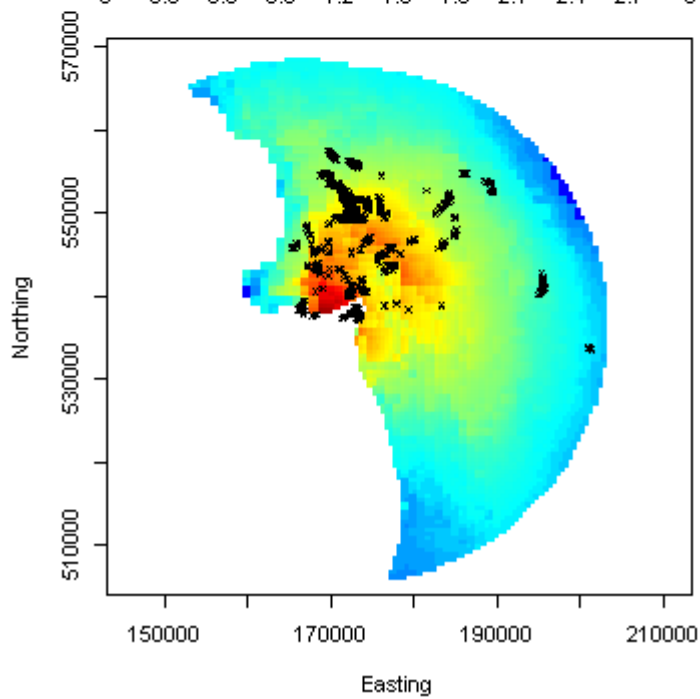
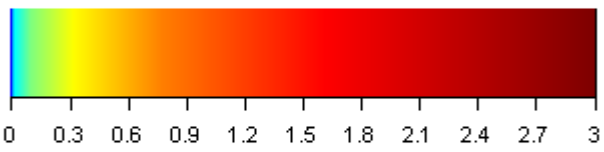
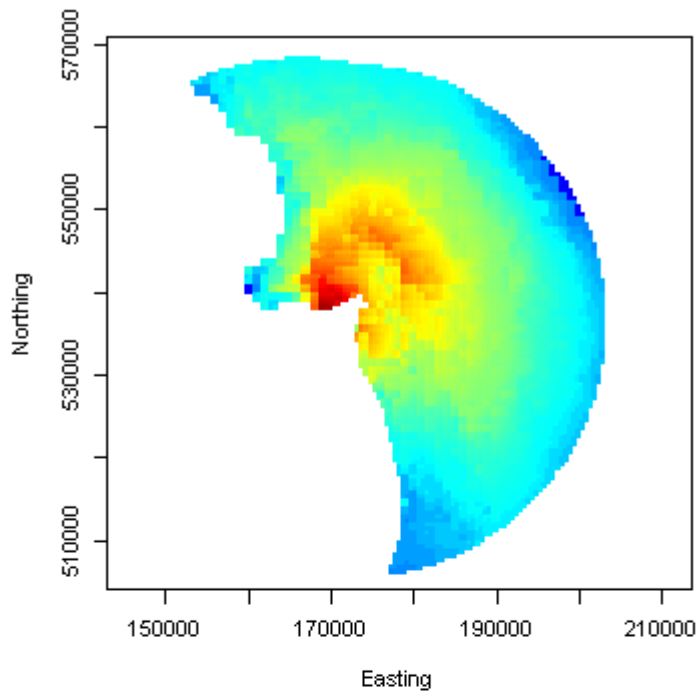
Arctic Tern, Farnes Colony – Usage



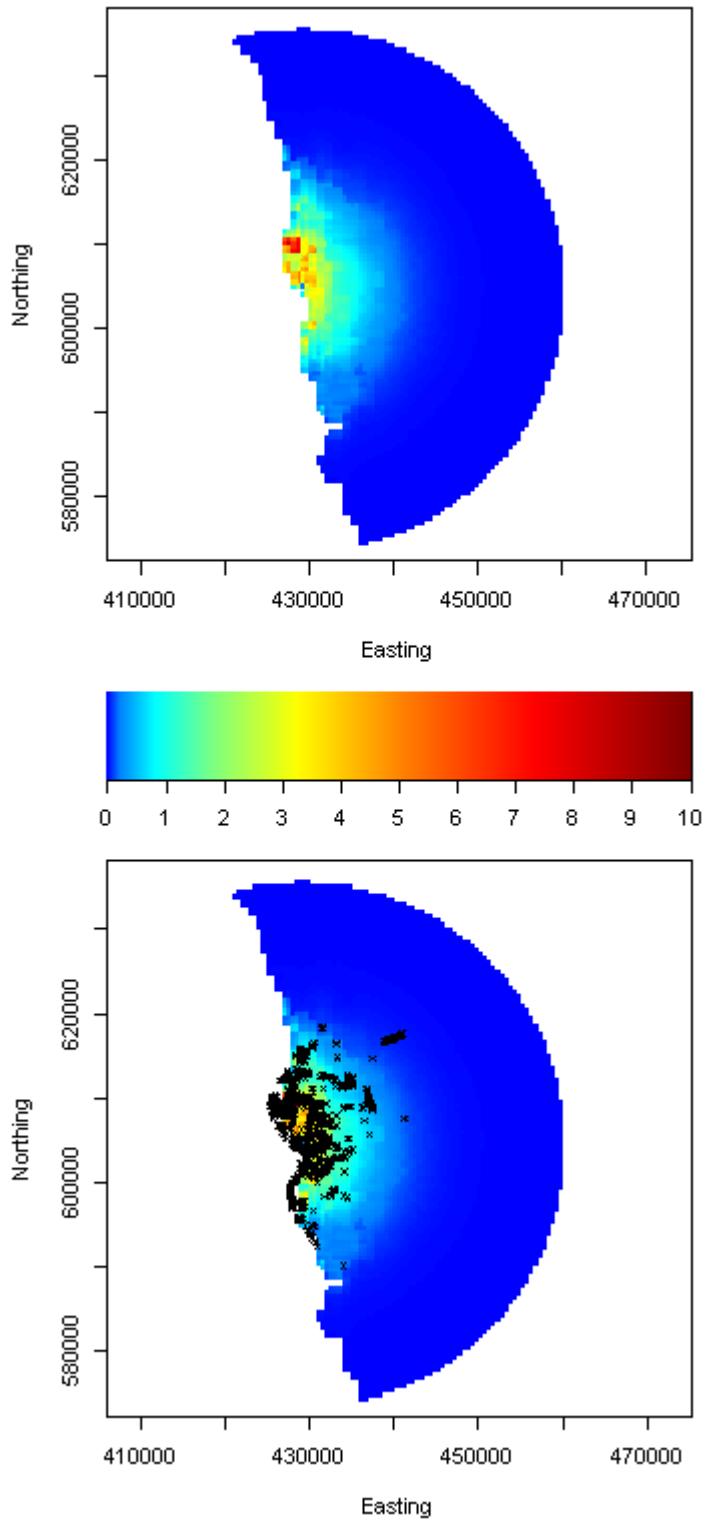
Arctic Tern, Outer Ards Colony – Preference



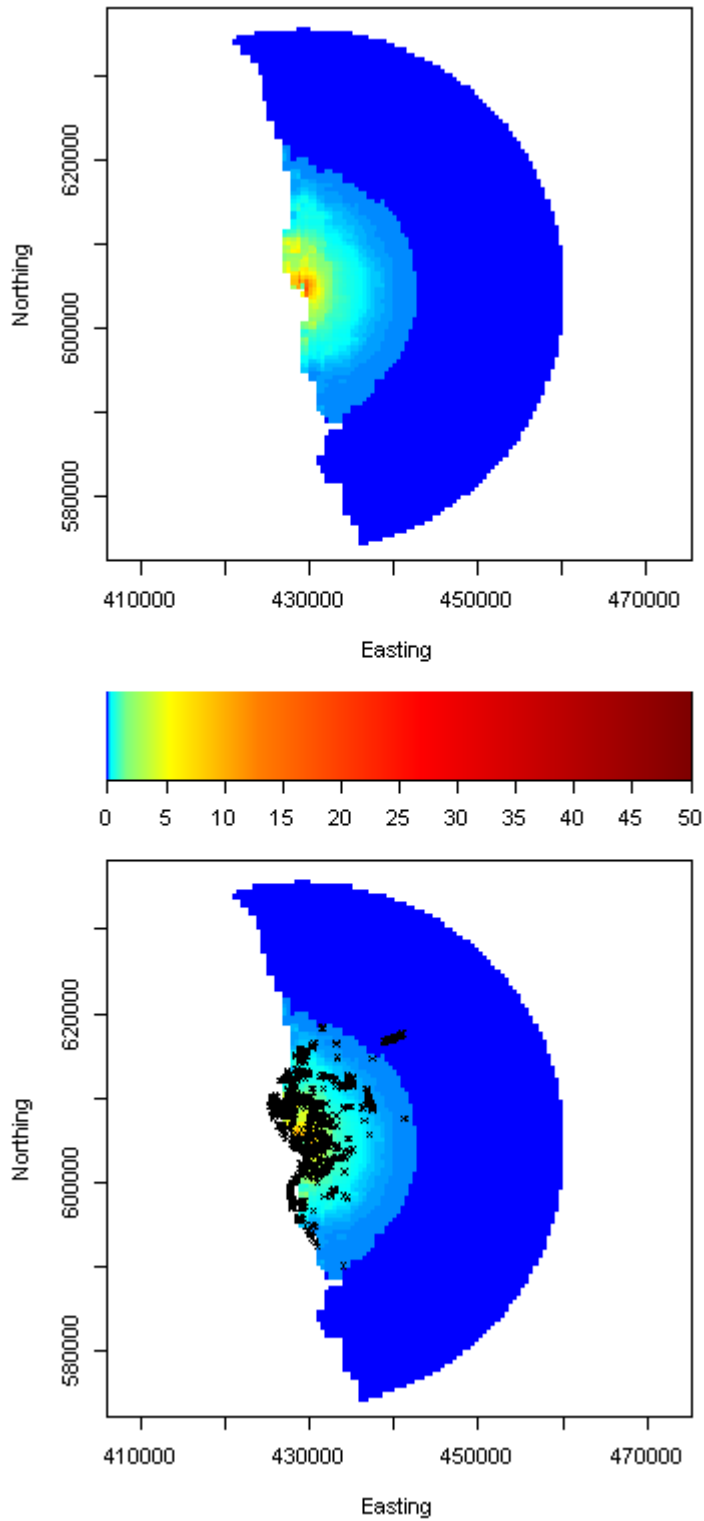
Arctic Tern, Outer Ards Colony – Usage



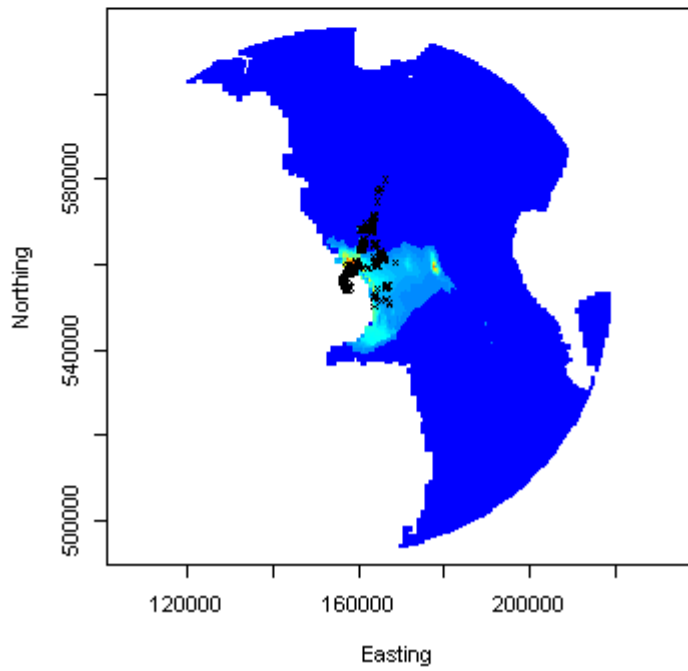
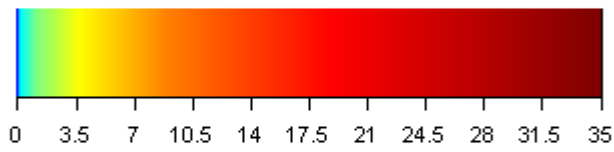
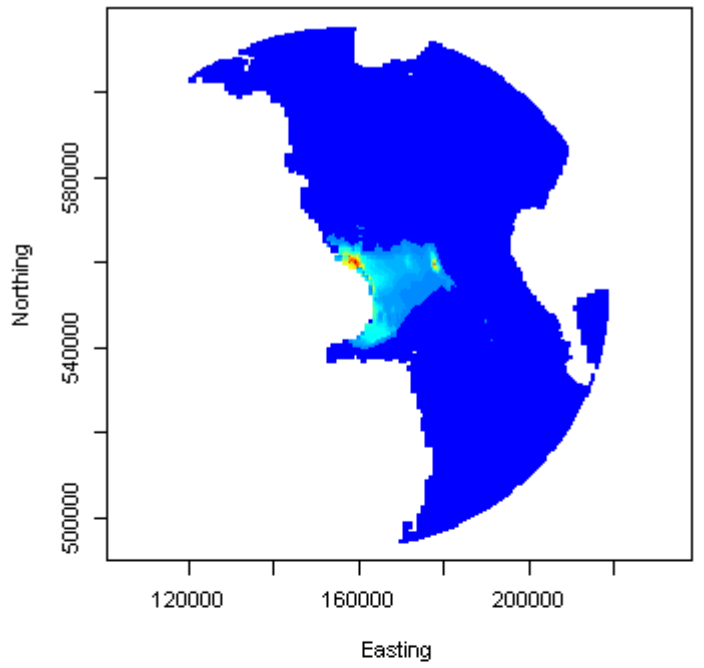
Common Tern, Coquet Colony – Preference



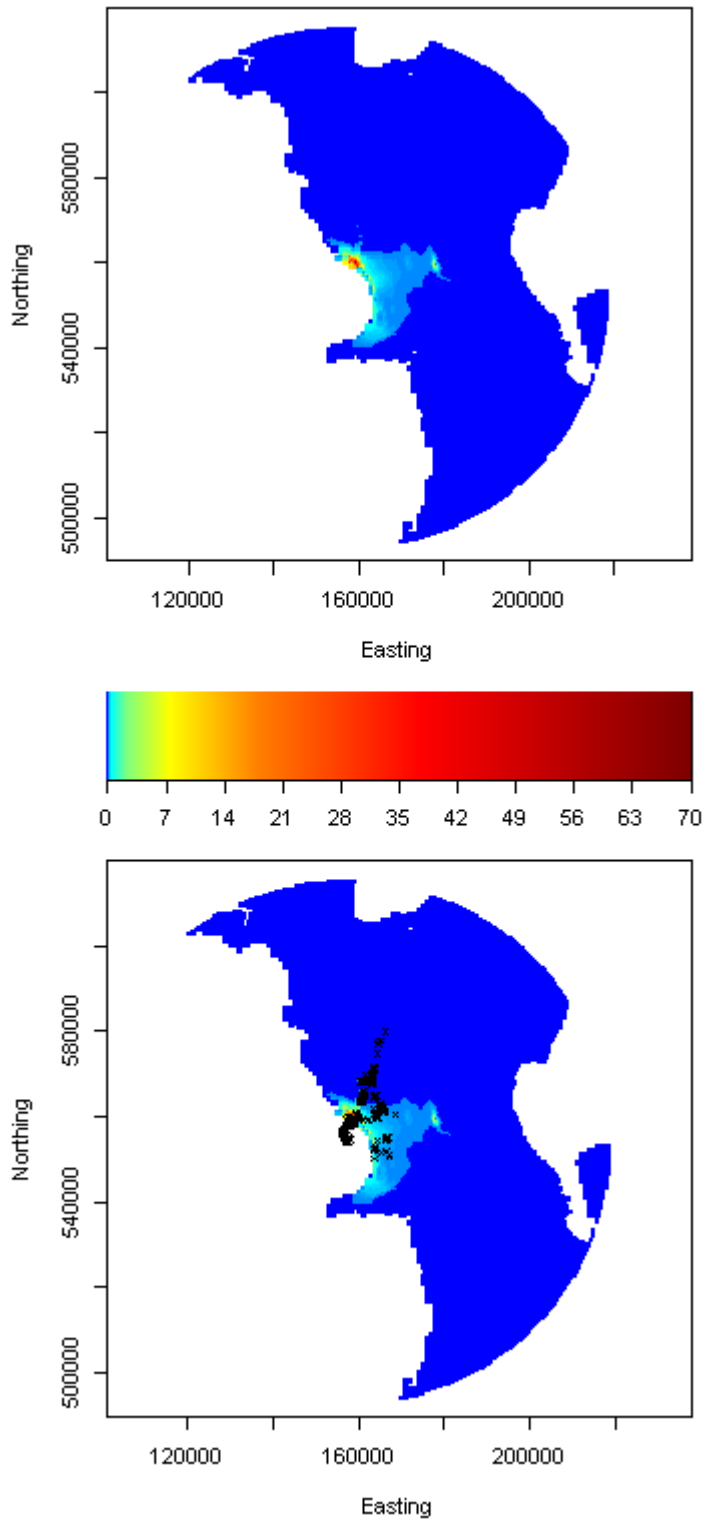
Common Tern, Coquet Colony – Usage



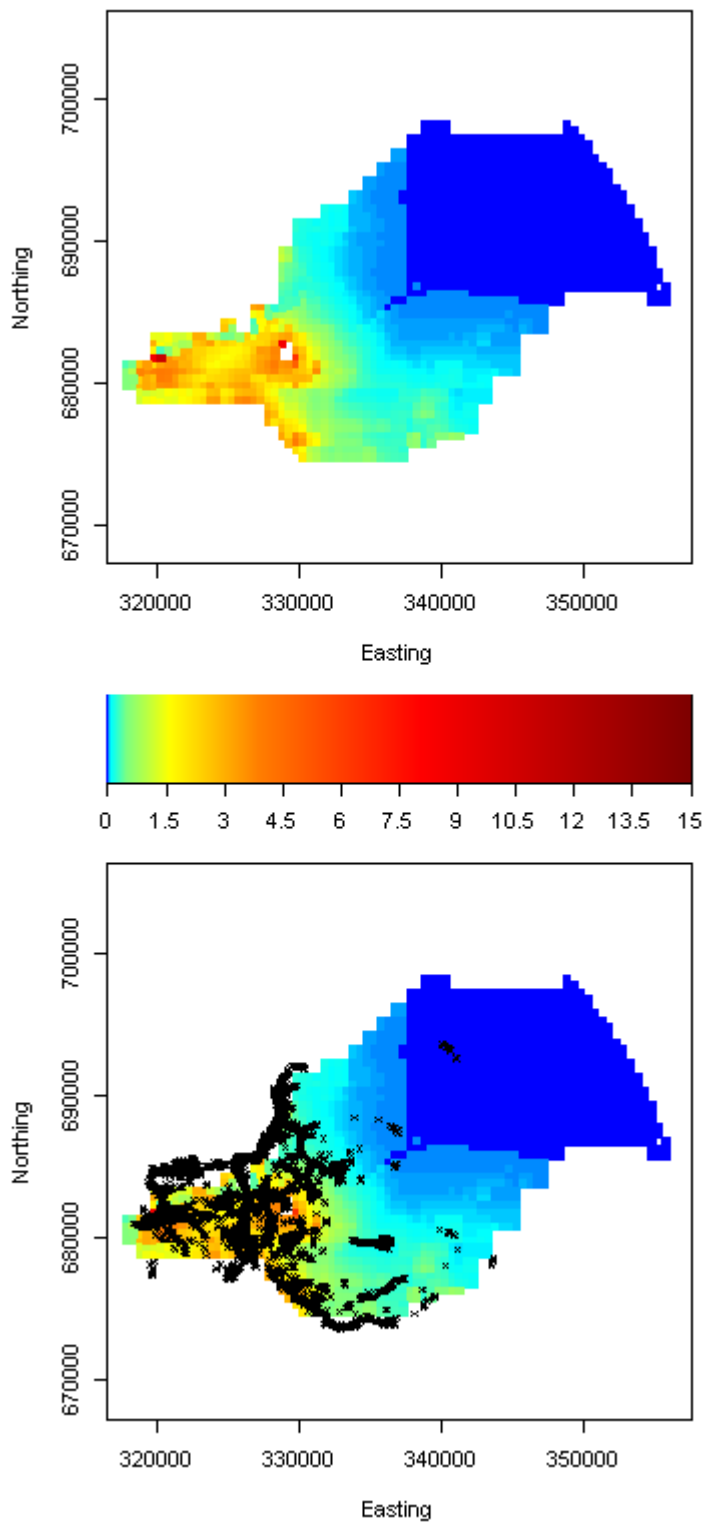
Common Tern, Larne Lough Colony – Preference



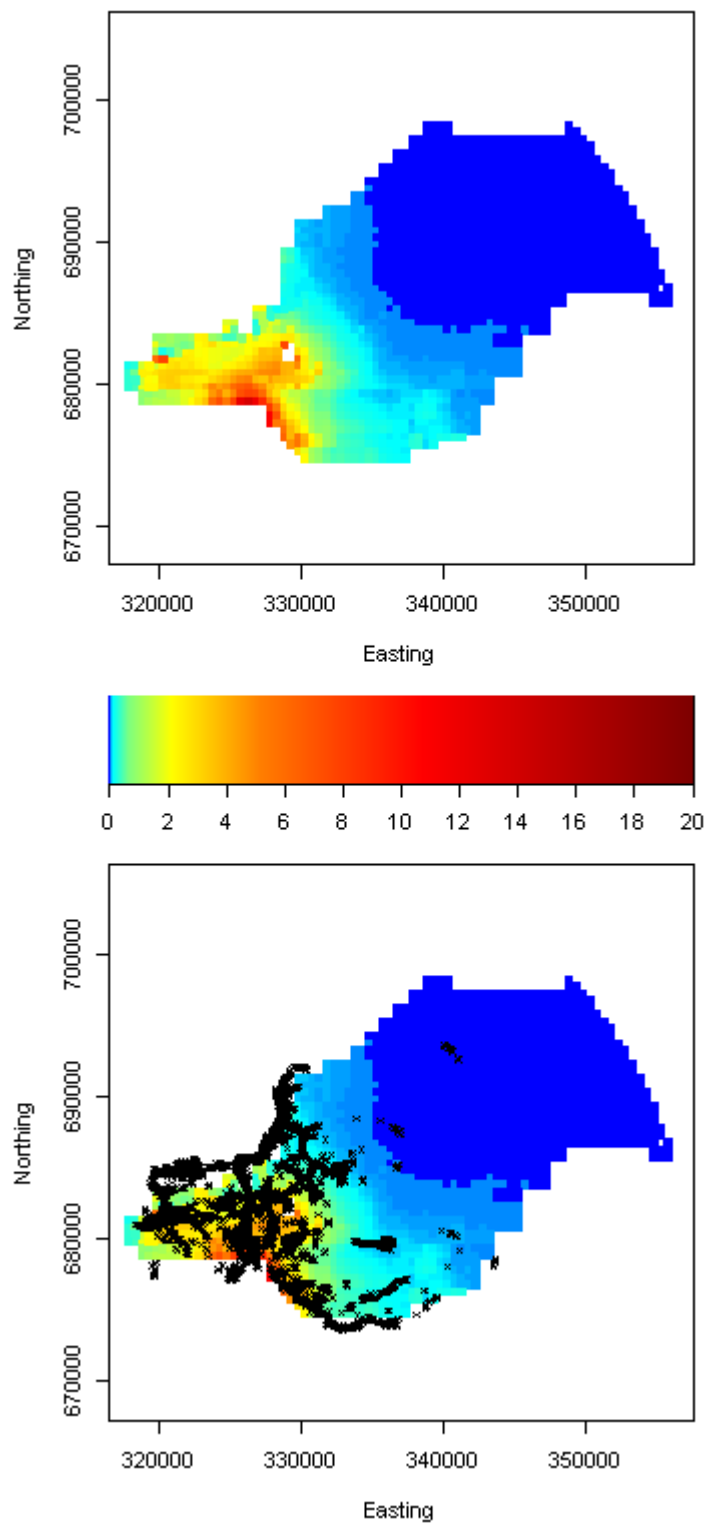
Common Tern, Larne Lough Colony – Usage



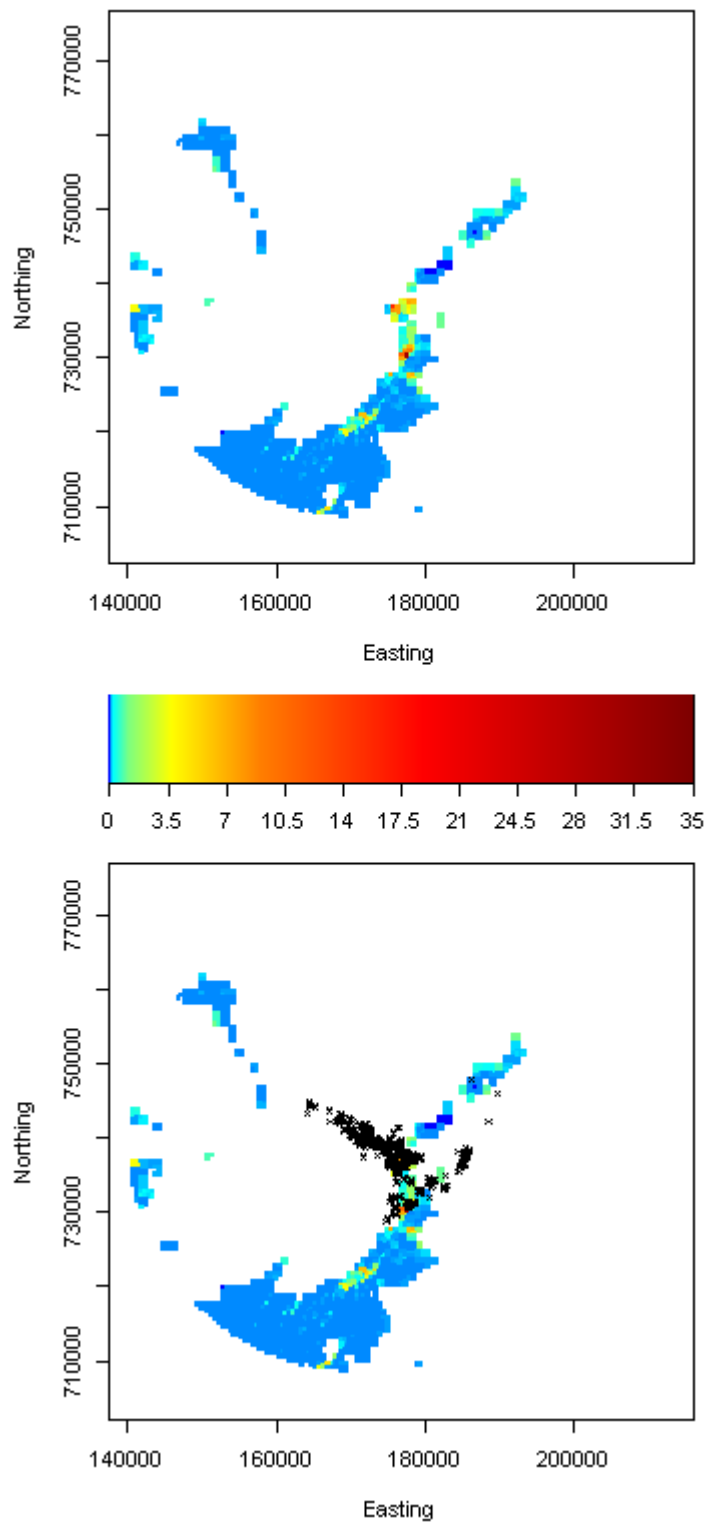
Common Tern, Leith Colony – Preference



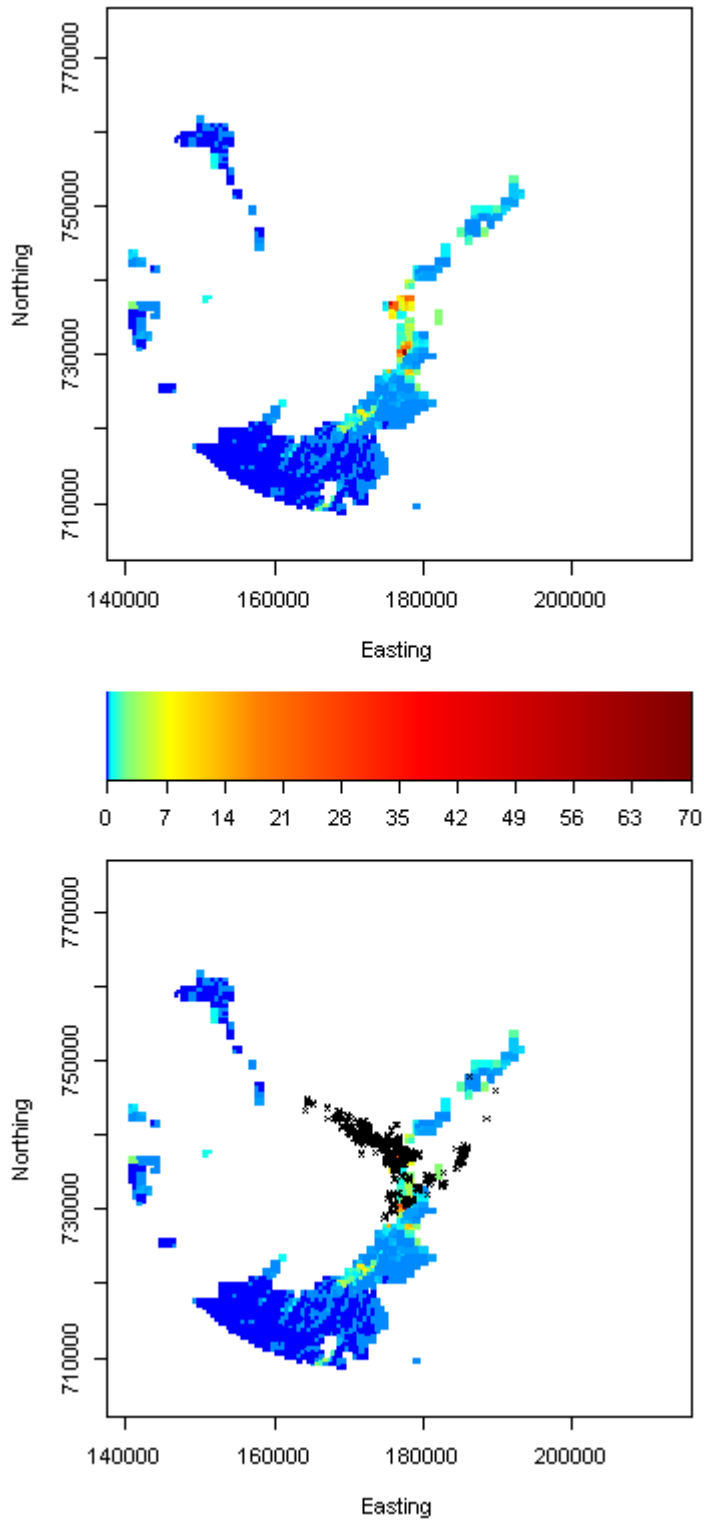
Common Tern, Leith Colony – Usage



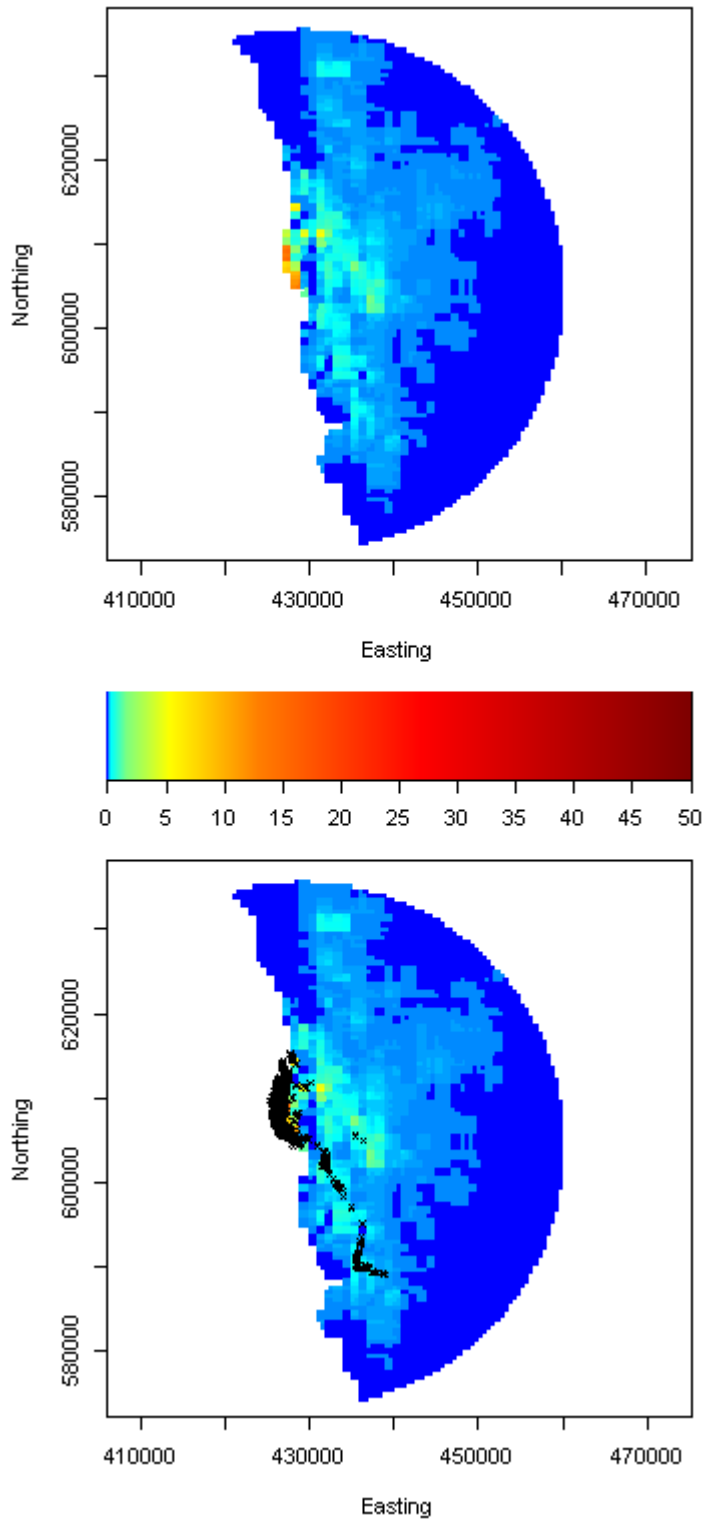
Common Tern, Mull Colony – Preference



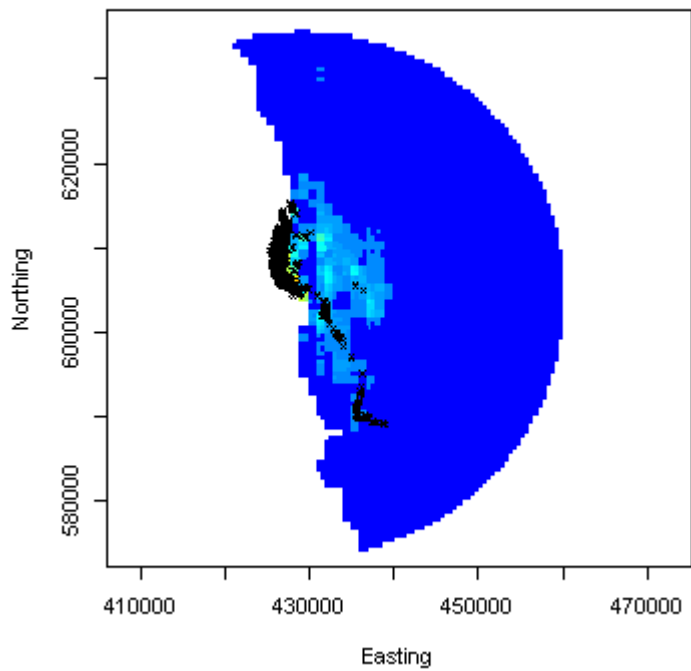
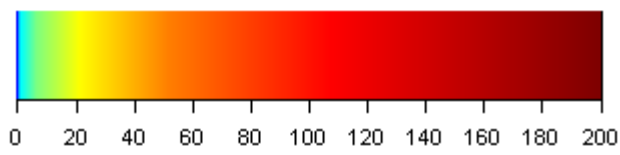
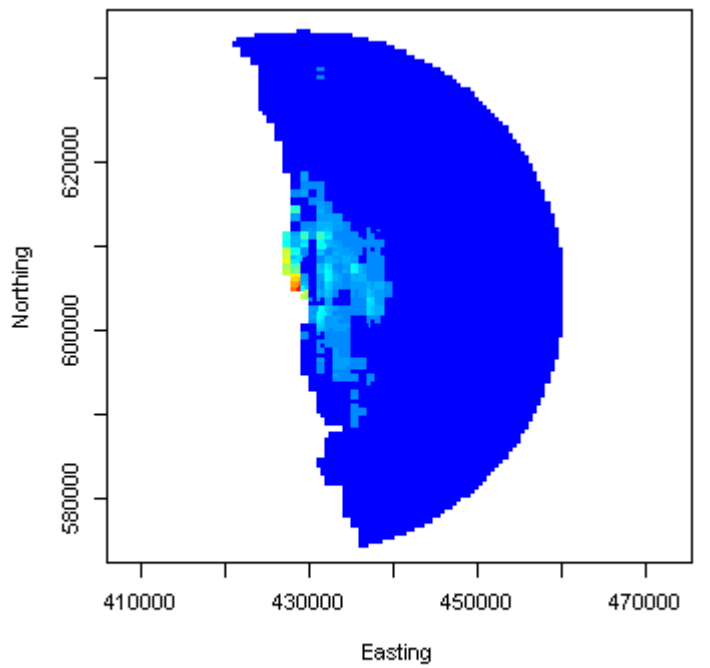
Common Tern, Mull Colony – Usage



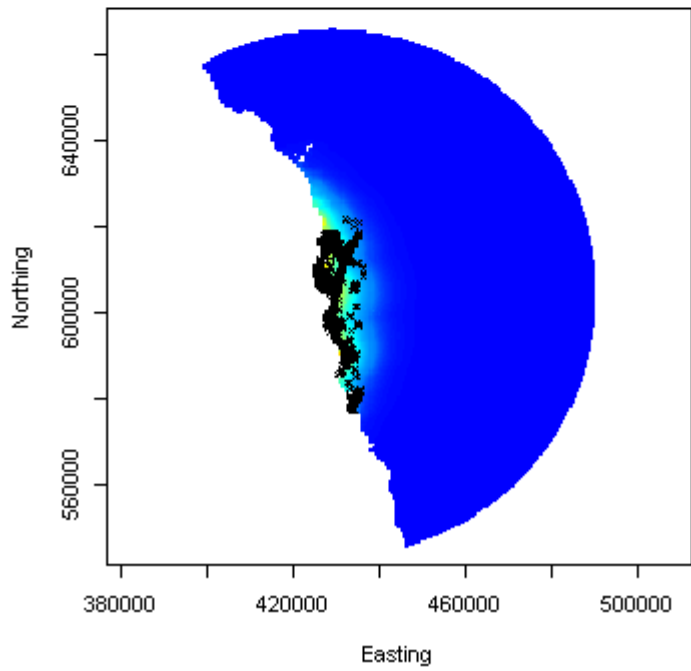
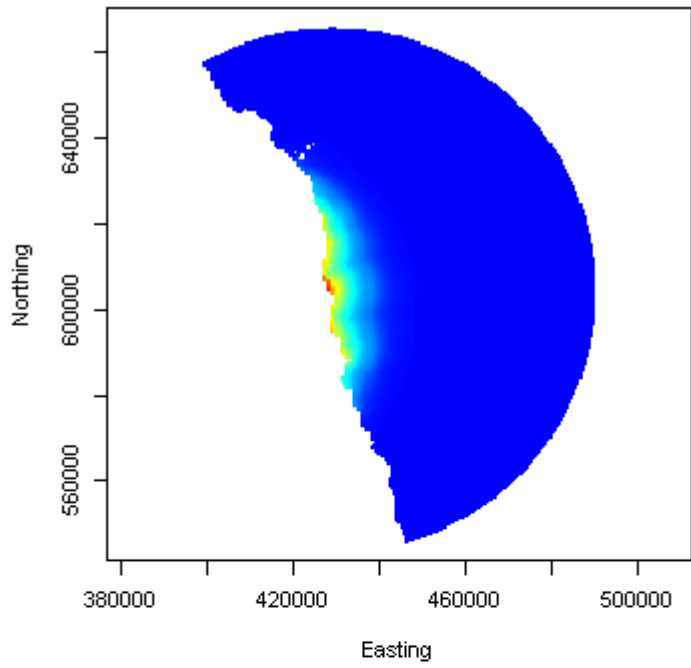
Roseate Tern, Coquet Colony – Preference



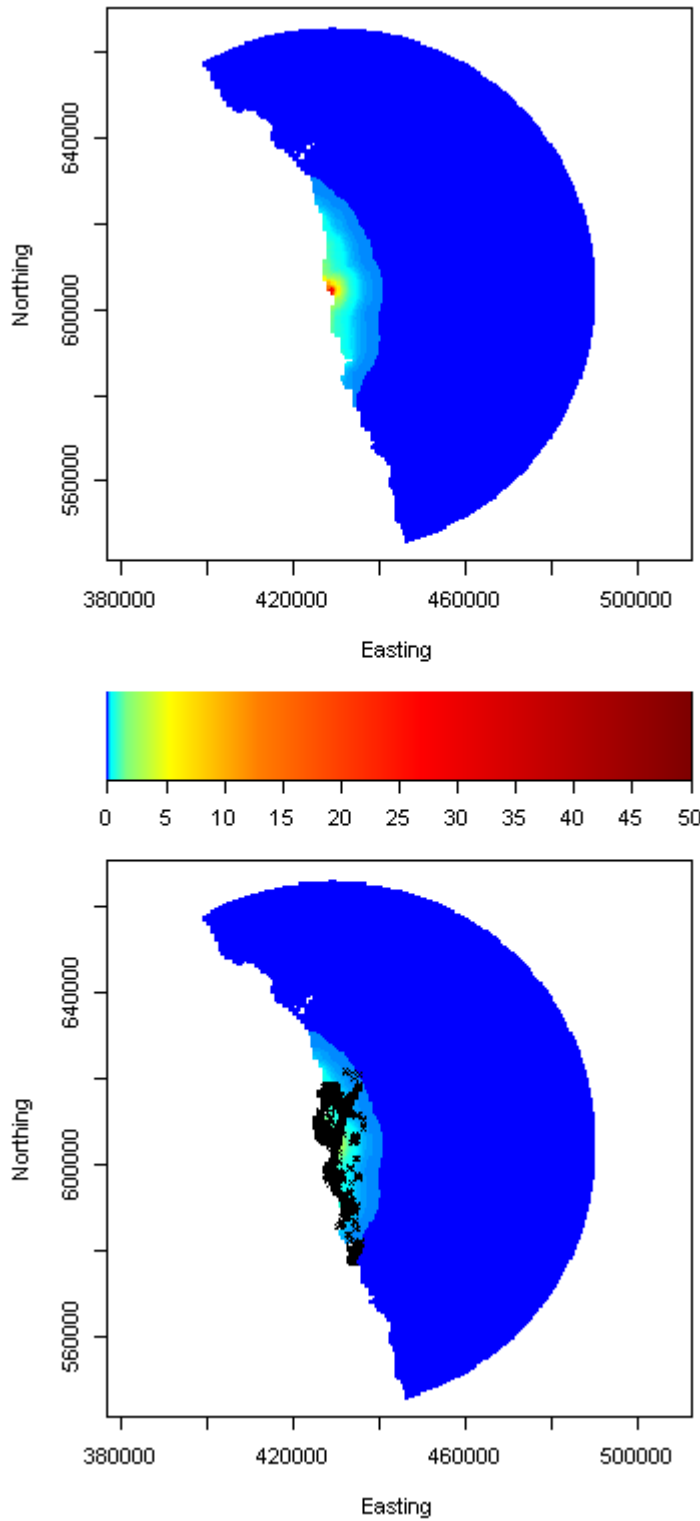
Roseate Tern, Coquet Colony – Usage



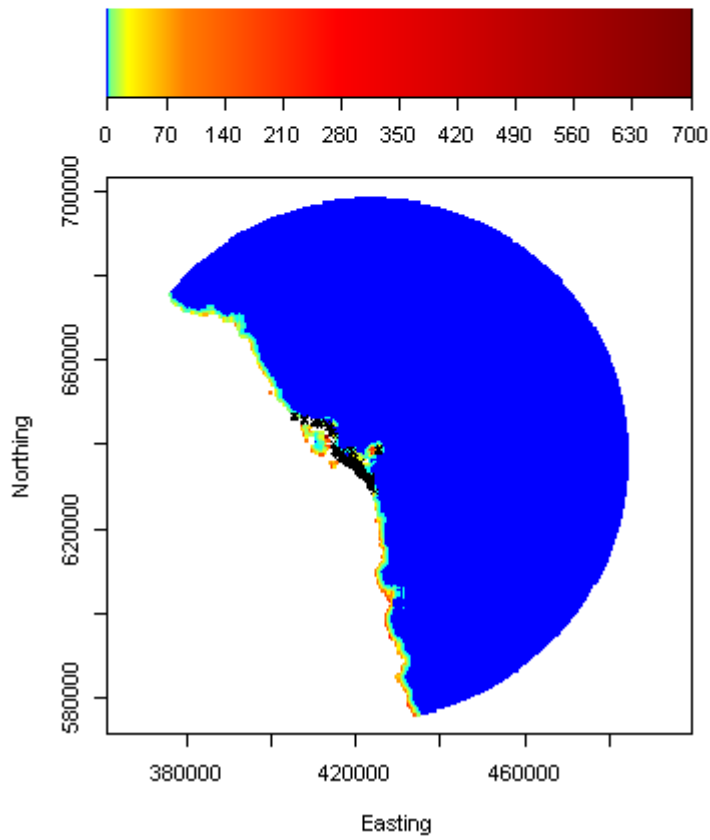
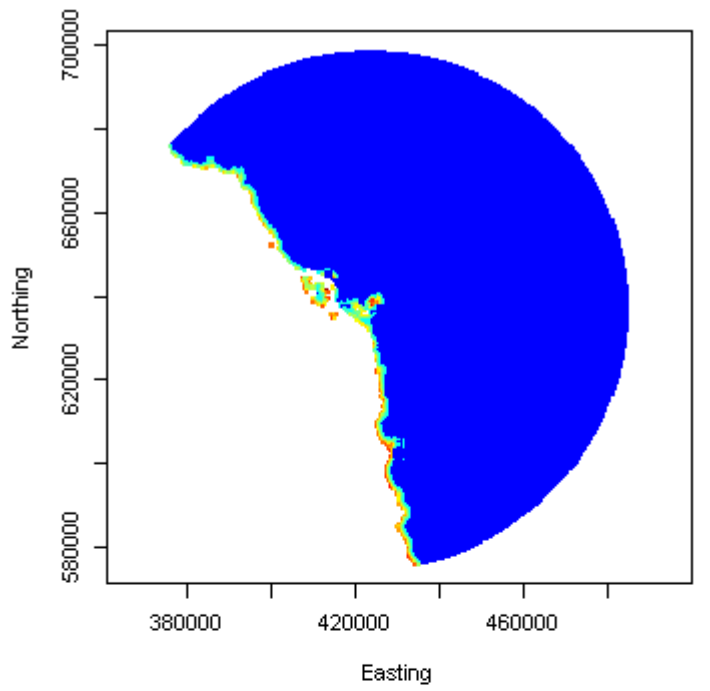
Sandwich Tern, Coquet Colony – Preference



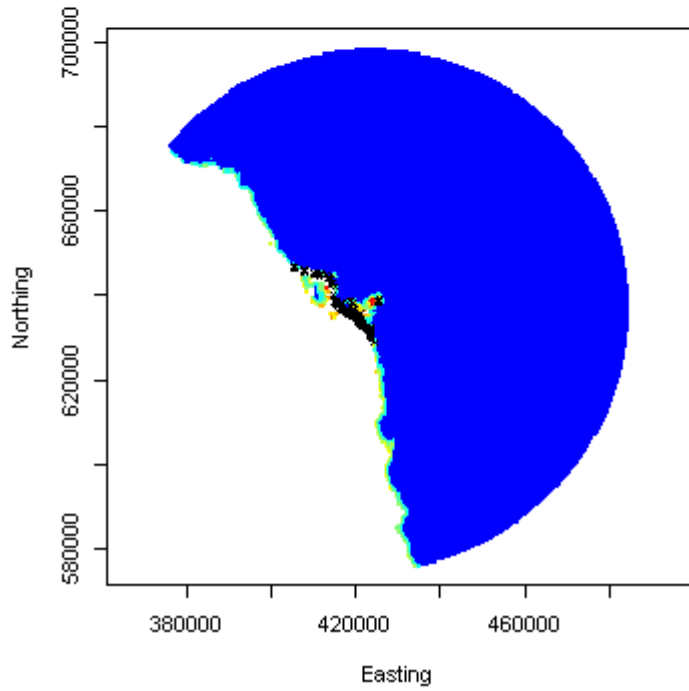
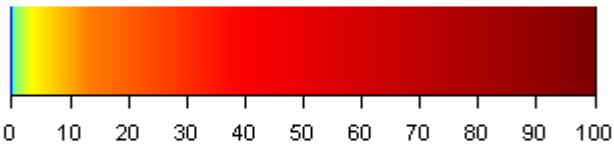
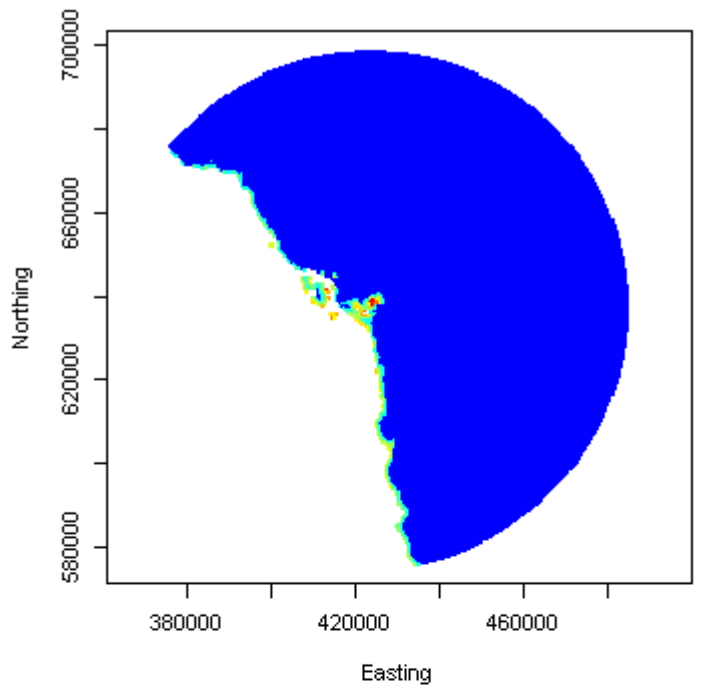
Sandwich Tern, Coquet Colony – Usage



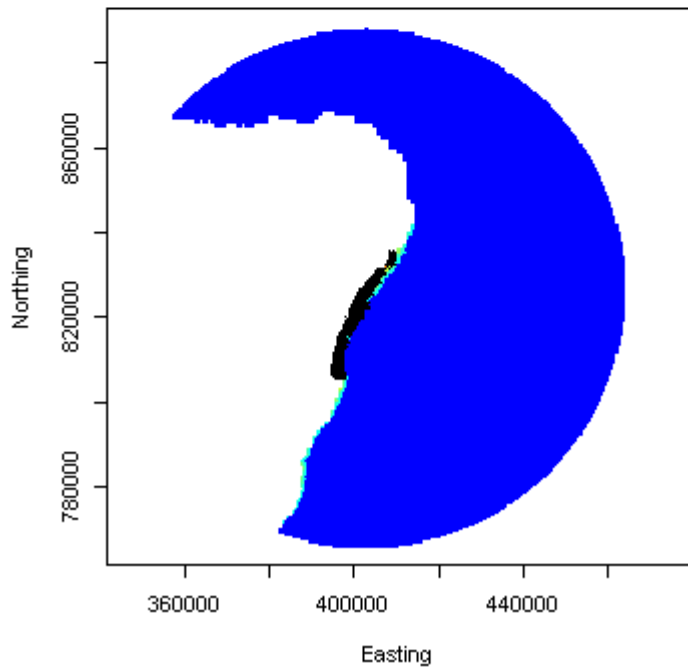
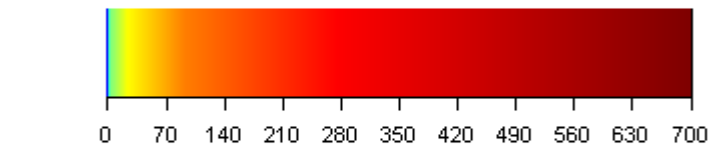
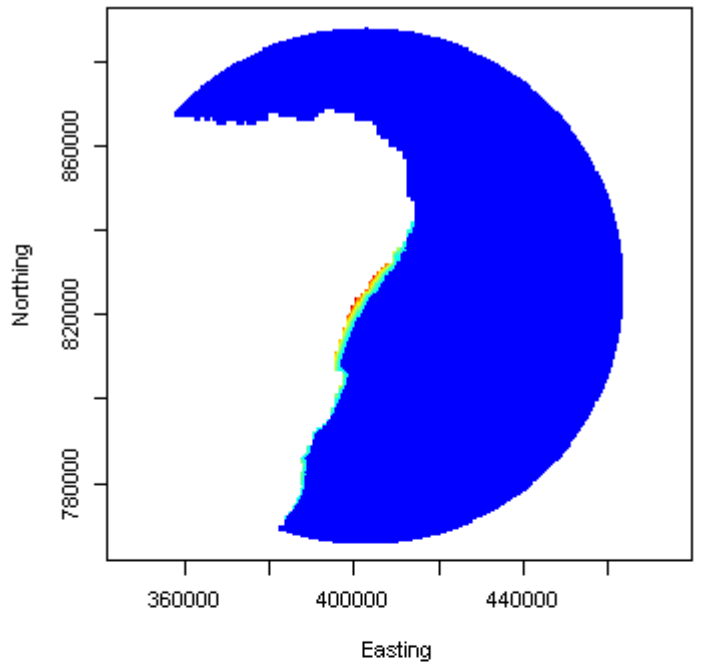
Sandwich Tern, Farnes Colony – Preference



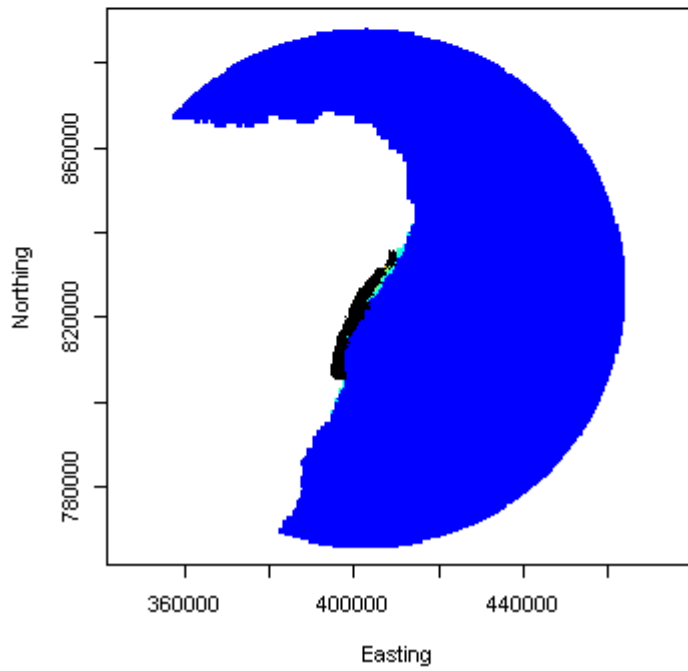
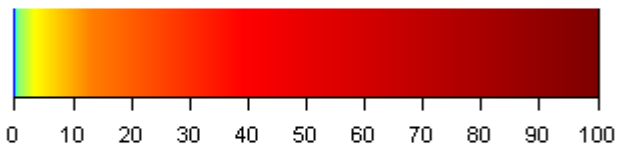
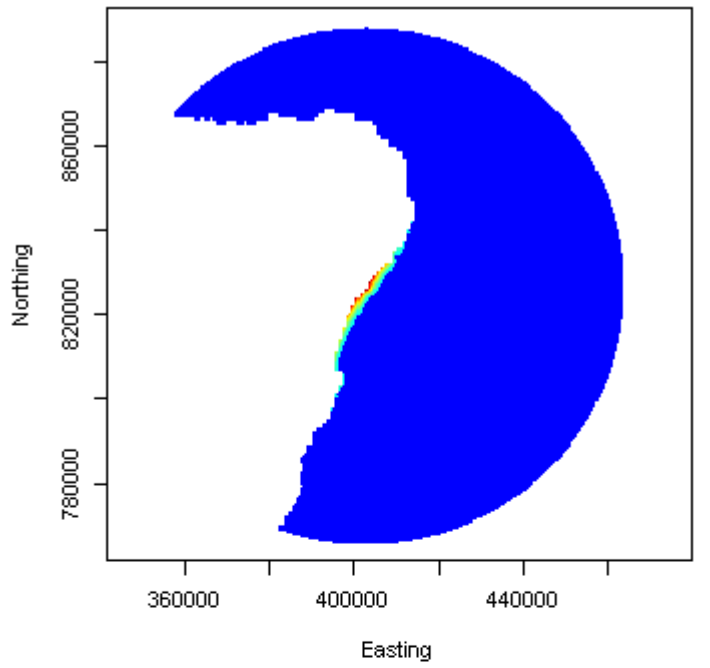
Sandwich Tern, Farnes Colony – Usage



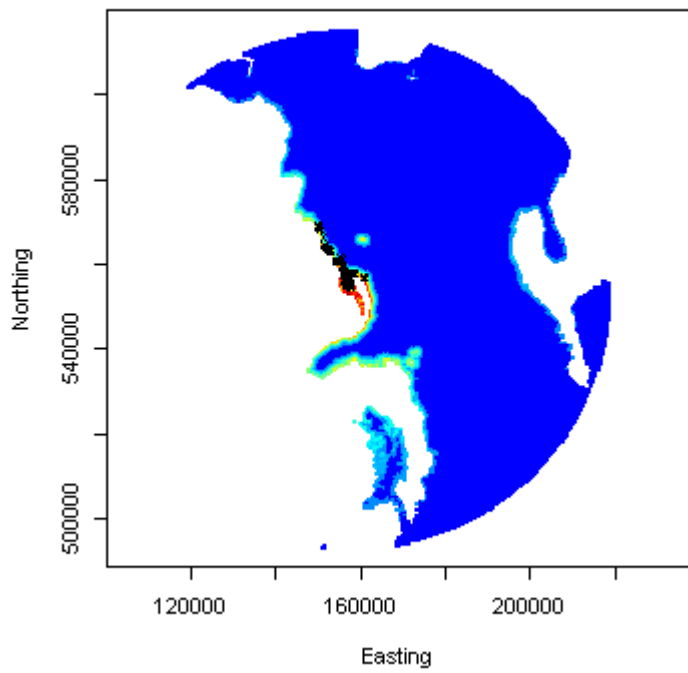
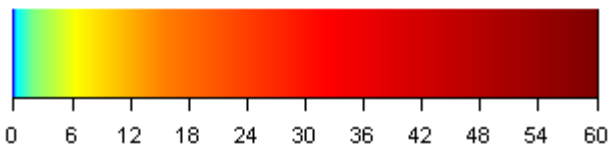
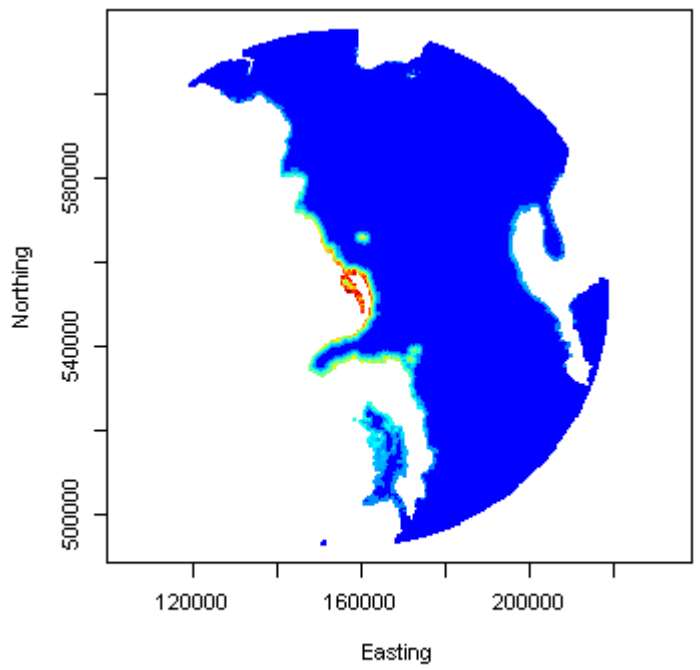
Sandwich Tern, Forvie Colony – Preference



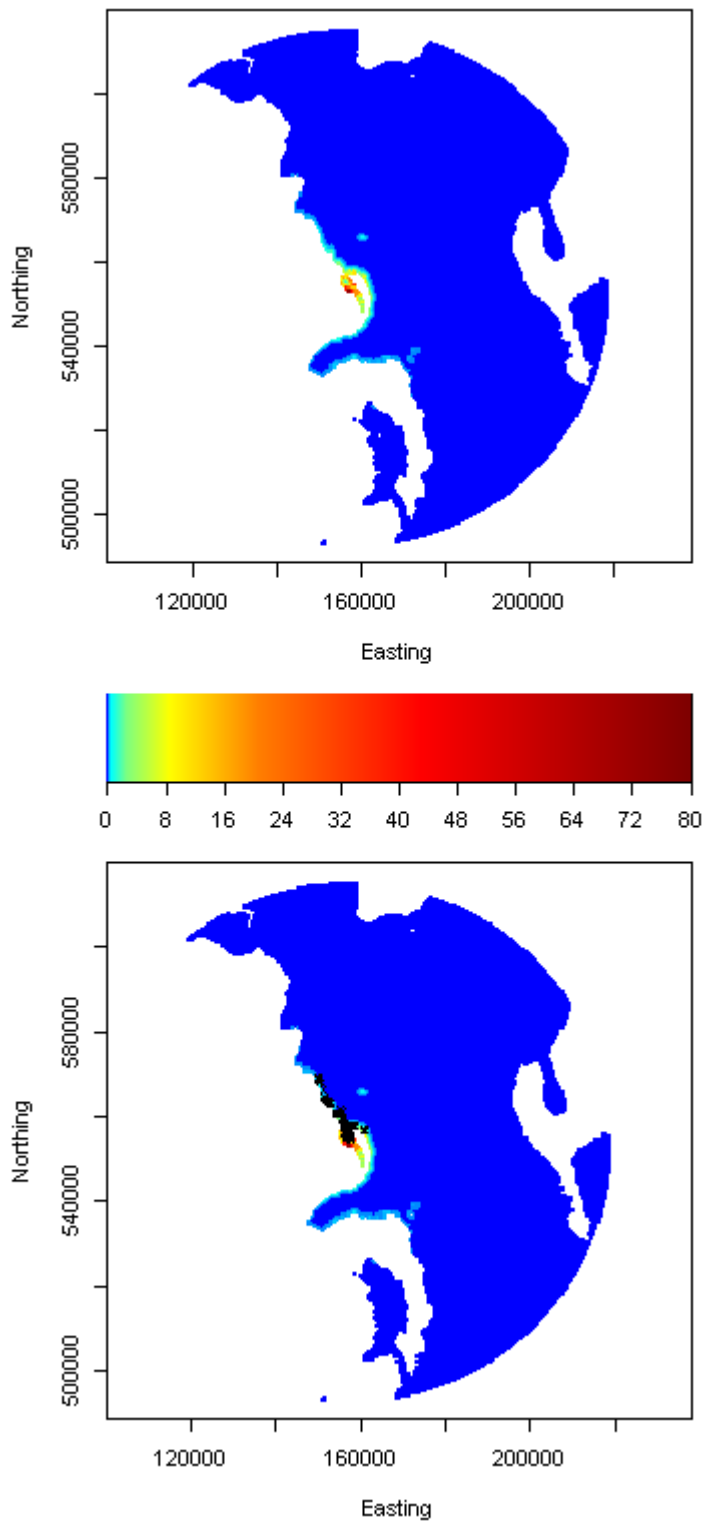
Sandwich Tern, Forvie Colony – Usage



Sandwich Tern, Larne Lough Colony – Preference



Sandwich Tern, Larne Lough Colony – Usage



References

- Aarts, G., MacKenzie, M., McConnell, B., Fedak, M. and Matthiopoulos, J. (2008) Estimating space-use and habitat preference from wildlife telemetry data. *Ecography*, **31**, 140-160.
- Aarts, G., Fieberg, J. and Matthiopoulos, J. (2012). Comparative interpretation of count, presence–absence and point methods for species distribution models. *Methods in Ecology and Evolution*, **3**, 177-187.
- Beale, C.M., Lennon, J.J., Yearsley, J.M., Brewer, M.J. and Elston, D.A. (2010) Regression analysis of spatial data. *Ecology Letters*, **13**, 246-264.
- Boyce, M.S. and McDonald, L.L. (1999) Relating populations to habitats using resource selection functions. *Trends in Ecology and Evolution*, **14**, 268-272.
- INLA (2012) R-Package, <http://www.r-inla.org/> .
- Keating, K. A. and Cherry, S. (2004) Use and interpretation of logistic regression in habitat selection studies. *Journal of Wildlife Management*, **68**, 774-789.
- R Development Core Team (2012) *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, <http://www.R-project.org/> .
- Shmueli, G. (2010) To explain or to predict? *Statistical Science*, **25**, 289-310.
- Wakefield, E.D., Phillips, R.A., Trathan, P.N., Arata, J., Gales, R., Huin, N., Robertson, G., Waugh, S.M., Weimerskirch, H. and Matthiopoulos, J. (2011) Habitat preference, accessibility, and competition limit the global distribution of breeding Black-browed Albatrosses. *Ecological Monographs*, **81**, 141-167.
- Warton, D.I. and Shepherd, L.C. (2010) Poisson point process models solve the “pseudo-absence problem” for presence-only data in ecology. *Annals of Applied Statistics*, **4**, 1383-1402.
- Wood, S.N. (2011) Fast stable restricted maximum likelihood and marginal likelihood estimation of semiparametric generalized linear models. *Journal of the Royal Statistical Society B*, **73**, 3-36.