

Avoiding extrapolation bias when using statistical models to make ecological prediction

PAUL B. CONN^{1,2} AND DEVIN S. JOHNSON¹, PETER L. BOVENG¹

¹*National Marine Mammal Laboratory, NOAA, National Marine Fisheries Service, Alaska*

Fisheries Science Center, 7600 Sand Point Way NE, Seattle, WA 98115 USA

¹ *Abstract.* We'll do this later

² *Key words:* *Abundance, Extrapolation, Forecasting, General additive models,*

³ *Generalized linear models, Independent Variable Hull, Leverage, Occupancy, Spatial*

⁴ *regression*

INTRODUCTION

⁵ In ecology and conservation, a common goal is to make predictions about an

⁶ unsampled random variable given a limited sample from the target population. For

⁷ instance, given a model (\mathcal{M}), estimated parameters ($\hat{\boldsymbol{\theta}}$), and a covariate vector \mathbf{x}_i , we often

⁸ desire to predict a new observation y_{new} at i . For instance, we might use a generalized

⁹ linear model (McCullagh and Nelder, 1989) or one of its extensions to predict species

¹⁰ density or occurrence in a new location, or to predict the future trend of a population.

¹¹ Early in their training, ecologists and statisticians are warned against extrapolating

¹² statistical relationships past the range of observed data. This caution is easily interpreted

¹³ in the context of single-variable linear regression analysis; one should be cautious in using

¹⁴ the fitted relationship to make predictions at some new point y_{new} whenever $x_{new} < \min(\mathbf{x})$

¹⁵ or $x_{new} > \max(\mathbf{x})$. But what about more complicated situations where there are multiple

¹⁶ explanatory variables, or when one uses a spatial regression model to account for the

¹⁷ residual spatial autocorrelation that is inevitably present in patchy ecological data

²Email: paul.conn@noaa.gov

18 (Lichstein et al., 2002)? How reliable are spatially- or temporally-explicit predictions in
19 sophisticated models for animal abundance and occurrence?

20 Statisticians have long struggled with the conditions under which fitted regression
21 models are capable of making robust predictions at new combinations of explanatory
22 variables. The issue is sometimes considered more of a philosophical problem than a
23 statistical one, and has even been likened to soothsaying (Ehrenberg and Bound, 1993). To
24 our mind, the reliability of predictions from statistical models is likely a function of several
25 factors, including (i) the intensity of sampling, (ii) spatial or temporal proximity of the
26 prediction location to locations where there are data, (iii) variability of the ecological
27 process, and (iv) the similarity of explanatory covariates in prediction locations when
28 compared to the ensemble of covariates for observed data locations.

29 Our aim in this paper is to investigate extrapolation bias in the generalized linear
30 model and its extensions, including generalized additive models (GAMs; Hastie and
31 Tibshirani, 1999; Wood, 2006) and spatial, temporal, or spatio-temporal regression models
32 (STRMs). In particular, we exploit some of the same ideas used in multiple linear
33 regression regarding leverage and outliers (Cook, 1979) to operationally define
34 “extrapolation” as making predictions that occur outside of a generalized independent
35 variable hull (gIVH) of observed data points. Application of the gIVH can provide
36 intuition regarding the reliability of predictions in unobserved locations, and can aid in
37 model construction. Also, since the gIVH can be constructed solely with knowledge of
38 sampled locations and explanatory covariates (i.e., it does not necessarily require any
39 observed response variables), it can also be used to help guide survey design. We illustrate
40 use of the gIVH on a simulated occupancy dataset, on a species distribution model (SDM)
41 for ribbon seals in the eastern Bering Sea, and on a population trend model for Steller Sea

42 Lions (*Phoca fasciata*).

43 THE GENERALIZED INDEPENDENT VARIABLE HULL (GIVH)

44 Extrapolation is often distinguished from interpolation. In a prediction context, we might
45 define (admittedly quite imprecisely) that extrapolation consists of making predictions that
46 are “outside the range of observed data” while interpolation consists of making predictions
47 “inside the range of observed data.” But what exactly do we mean by “outside the range of
48 observed data”? Predictions outside the range of observed covariates? Predictions for
49 locations that are so far from places where data are gathered that we are skeptical that the
50 estimated statistical relationship still holds? To help guide our choice of an operational
51 definition, we turn to early work on outlier detection in simple linear regression analysis.

52 In the context of outlier detection, Cook (1979) defined an independent variable hull
53 (IVH) as the smallest convex set containing all design points of a full-rank linear regression
54 model. Linear regression models are often written in matrix form, i.e.

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

55
56 where \mathbf{Y} give observed data, \mathbf{X} is a so-called design matrix that includes explanatory
57 variables (see e.g. Draper and Smith, 1966), and $\boldsymbol{\epsilon}$ represent normally distributed residuals
58 (here and throughout the paper, bold symbols will be used to denote vectors and
59 matrices). Under this formulation, the IVH is defined relative to the hat matrix,
60 $\mathbf{V}_{LR} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ (where the subscript “LR” denotes linear regression). Letting v
61 denote the maximum diagonal element of \mathbf{V}_{LR} , one can examine whether a new design

62 point, \mathbf{x}_0 is within the IVH. In particular, \mathbf{x}_0 is within the IVH whenever

63
$$\mathbf{x}'_0(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_0 \leq v. \quad (1)$$

64 Cook (1979) used this concept to identify influential observations and possible outliers,

65 arguing that design points near the edge of the IVH are deserving of special attention.

66 Similarly, points outside the IVH should be interpreted with caution.

67 We simulated two sets of design data to help illustrate application of the IVH (Fig. 2).

68 In simple linear regression with one predictor variable, predictions on a hypothetical
69 response variable obtained at covariate values below the lowest observed value or above the
70 highest observed value are primarily outside of the IVH. We suspect this result conforms to
71 most ecologists intuition about what constitutes “extrapolating past the observed data.”

72 However, fitting a quadratic model exhibits more nuance; if there is a large gap between
73 design points, it is entirely possible that intermediate covariate values will also be outside
74 of the IVH and thus more likely to result in problematic predictions. Fitting a model with
75 two covariates and both linear and quadratic effects, the shape of the IVH is somewhat

76 more irregular. These simple examples highlight the sometimes counterintuitive nature of
77 the predictive inference problem, a problem that can only become worse as models with
78 more dimensions are contemplated (including those with temporal or spatial structure).

79 Fortunately, the ideas behind the IVH provide potential way forward.

80 Cook’s (1979) formulation for the IVH is particular to linear regression analysis, which
81 assumes iid normally distributed error. Thus, it is not directly applicable to generalized
82 models, such as those including spatial random effects. Further, the hat matrix is not
83 necessarily well defined for models with more general spatial structure. However, since the

84 hat matrix is proportional to prediction variance, Cook (1979) notes that design points
 85 with maximum prediction variance will be located on the boundary of the IVH. Working
 86 on the linear predictor scale, we therefore define a generalized independent variable hull
 87 (gIVH) as the set of all points \mathbf{x} (note that \mathbf{x} can include both observed and unobserved
 88 design points) such that

$$89 \quad \text{var}(\boldsymbol{\mu}_x | \mathbf{x}) \leq \max[\text{var}(\boldsymbol{\mu}_X | \mathbf{X})], \quad (2)$$

90 where $\boldsymbol{\mu}_x$ correspond to predictions at \mathbf{x} , and \mathbf{X} and $\boldsymbol{\mu}_X$ denote observed design points and
 91 predictions at X , respectively.

92 Generalizations of the linear model are often written in the form

$$93 \quad Y_i \sim f_Y(g^{-1}(\mu_i)), \quad (3)$$

94 where f_Y denotes a probability density or mass function (e.g. Bernoulli, Poisson), g gives a
 95 link function, and μ_i is a predictor. For many such generalizations, it is possible to specify
 96 the μ_i as

$$97 \quad \boldsymbol{\mu} = \mathbf{X}_{aug} \boldsymbol{\beta}_{aug} + \boldsymbol{\epsilon}, \quad (4)$$

98 where the $\boldsymbol{\epsilon}$ represent Gaussian error, \mathbf{X}_{aug} denotes an augmented design matrix, and $\boldsymbol{\beta}_{aug}$
 99 denote an augmented vector of parameters. For instance, in a spatial model, $\boldsymbol{\beta}_{aug}$ might
 100 include both fixed effect parameters and spatial random effects in a reduced dimension
 101 subspace (see Appendix S1 for examples of how numerous types of models can be written
 102 in this form).

103 When models are specified as in Eq. 4, we can write prediction variance generically as

104 $\text{var}(\boldsymbol{\mu}_x | \mathbf{x}) = \mathbf{x}\text{var}(\hat{\boldsymbol{\beta}}_{\text{aug}})\mathbf{x}',$ (5)

105 where it is understood that the exact form of \mathbf{x} and $\text{var}(\hat{\boldsymbol{\beta}}_{\text{aug}})$ depends on the model chosen
106 (i.e., GLM, GAM, or STRM; Appendix S1). For GLMs, prediction variance is proportional
107 to the original hat matrix \mathbf{V}_{LR} (at least on the linear predictor scale), and Eq. 1 may be
108 used directly. For other models, we have

109 $\text{var}(\hat{\boldsymbol{\beta}}_{\text{aug}}) = \begin{bmatrix} \Sigma_\beta & \mathbf{0} \\ \mathbf{0} & \Sigma_\alpha \end{bmatrix},$ (6)

110 where we simply replace $\boldsymbol{\beta}_{\text{aug}}$ with the augmented vector that includes both regression and
111 spatial or smooth parameters (i.e. $\boldsymbol{\beta}$ and $\boldsymbol{\alpha}$; Appendix S1).

112 The exact form of $\text{var}(\hat{\boldsymbol{\beta}}_{\text{aug}})$ differs depending on the underlying model structure and
113 estimation procedure. In the following treatment, we shall focus on Bayesian analysis;
114 although it is not necessarily needed to fit GLMs and GAMs, it helps make STRMs more
115 computationally tractable and puts all of the different models into a common analysis
116 framework. This approach requires specifying priors for model parameters, and prior
117 parameters may appear in expressions for prediction variance. Judicious choices of priors
118 can at times limit the influence of priors on calculation of the gIVH. For example, if we
119 specify the prior on the vector of regression coefficients to be $[\boldsymbol{\beta}] = \text{MVN}(\mathbf{0}, (\tau_\beta X'X)^{-1})$
120 where τ_β is a fixed constant, a Bayesian implementation of a GLM still results in the gIVH
121 specified in Eq. 1. The ability to use Eq. 5 directly is quite useful, as we only need to know
122 the explanatory variables to be able to diagnose whether predictions lie in or out of the

123 gIVH. However, the other models considered here (GAMs and STRMs) typically require
 124 estimation of an additional precision parameter. In these cases, one solution is to impose
 125 prior distributions and integrate over precision parameters (call these $\boldsymbol{\tau}$) to obtain an
 126 expectation for $\text{var}(\hat{\boldsymbol{\beta}}_{aug})$. Using the law of the unconscious statistician, we have

$$127 \quad \text{E}(\text{var}(\hat{\boldsymbol{\beta}}_{aug})) = \int_{\boldsymbol{\tau}} \text{var}(\hat{\boldsymbol{\beta}}_{aug})[\boldsymbol{\tau}] d\boldsymbol{\tau}. \quad (7)$$

128 Here, and throughout the paper, we use the bracket notation to indicate a probability
 129 density function; e.g., $[\boldsymbol{\tau}]$ denotes the joint prior distribution for precision parameters $\boldsymbol{\tau}$. If
 130 observations have already been gathered, one could also consider integrating over the
 131 marginal posterior distribution using MCMC:

$$132 \quad \text{E}(\text{var}(\hat{\boldsymbol{\beta}}_{aug})|\mathbf{Y}) = \int_{\boldsymbol{\tau}} \text{var}(\hat{\boldsymbol{\beta}}_{aug})[\boldsymbol{\tau}|\mathbf{Y}] d\boldsymbol{\tau}. \quad (8)$$

133 We propose to use the gIVH in much the same manner as Cook (1979). In particular,
 134 we use the gIVH to differentiate whether spatial predictions are interpolations (predictive
 135 design points lying inside the gIVH) or extrapolations (predictive design points lying
 136 outside the gIVH). To our knowledge, this is the first time this definition has been used in
 137 the context of spatial prediction, but one we shall show can be quite helpful in diagnosing
 138 and mitigating extrapolation bias. The gIVH seems ideally situated to this task as it does
 139 not necessarily need to rely on gathered response data. Thus, one can examine whether or
 140 not prediction points lie within the IVH without ever collecting response data there.

141

EXAMPLES

142

Simulation study

143 We conducted simulation analyses to investigate whether the gIVH was useful in
 144 diagnosing prediction biases in occupancy and abundance analyses. In both cases, we
 145 abundance on a 30×30 grid, and generated as a function of four habitat covariates in
 146 addition to spatial process covariance (Appendix S2). The habitat covariates also

147

Ribbon seal SDM

148 As part of an international effort, researchers with the U.S. National Marine Fisheries
 149 Service conducted aerial surveys over the eastern Bering Sea in 2012 and 2013. Agency
 150 scientists used infrared video to detect seals that were on ice, and simultaneous automated
 151 digital photographs provided information on species identity. Here, we use spatially
 152 referenced count data from photographed ribbon seals, *Phoca fasciata* (Fig. 1) on a subset
 153 of 10 flights flown over the Bering Sea in April 2012. These flights were limited to a one
 154 week period so that both sea ice conditions and seal distributions could be assumed to be
 155 static.

156 Our objective with this dataset will be to model seal counts on transects through 25km
 157 by 25km grid cells as a function of habitat covariates and possible spatial autocorrelation.
 158 Estimates of apparent abundance can then be obtained by summing predictions across grid
 159 cells. Figure 3 shows the transects flown and the number of ribbon seals encountered in
 160 each cell, and Figure 5 show explanatory covariates gathered to help predict ribbon seal
 161 abundance. These data are described in fuller detail by Conn et al. (Accepted), who
 162 extend the modeling framework of STRMs to account for incomplete detection and species

₁₆₃ misidentification errors (see e.g. Conn et al., Accepted). Since our focus in this paper is on
₁₆₄ illustrating spatial modeling concepts, we devote our efforts to the comparably easier
₁₆₅ problem of estimating apparent abundance (i.e., uncorrected for vagaries of the detection
₁₆₆ process).

₁₆₇ Inspection of ribbon seal data (Fig. 3) immediately reveals a potential issue with
₁₆₈ spatial prediction: abundance of ribbon seals appears to be maximized in the southern
₁₆₉ and/or southeast quadrant of the surveyed area. Predicting abundance in areas further
₁₇₀ south and east may thus prove problematic. To illustrate, let Y_i denote the ribbon seal
₁₇₁ count (Y_i) obtained in sampled grid cell i . Suppose that counts arise according to a
₁₇₂ log-Gaussian Cox process, such that

$$Y_i \sim \text{Poisson}(\lambda_i) \text{ and} \quad (9)$$

$$\log(\lambda_i) = \log(P_i) + \mathbf{X}_i\boldsymbol{\beta} + \eta_i + \epsilon_i,$$

₁₇₃ where P_i gives the proportion of area photographed in grid cell i (recall also that \mathbf{X}_i
₁₇₄ denotes a vector covariates for cell i , $\boldsymbol{\beta}$ are regression coefficients, η_i represents a spatially
₁₇₅ autocorrelated random effect, and ϵ_i is normally distributed *iid* error).

₁₇₆ We could fit any number of predictive models to these data, but we start with a simple
₁₇₇ generalized linear model where we ignore the spatial random effect, η_i , and use the full
₁₇₈ suite of predictor covariates (Fig. 5) to fit Eq. 9 to our data. In particular, we fit a model
₁₇₉ with linear effects of all predictor variables, and with an additional quadratic term for ice
₁₈₀ concentration (seal density is often maximized at an intermediate value of ice
₁₈₁ concentration; see Ver Hoef et al., 2013; Conn et al., Accepted). To enable comparison
₁₈₂ with more complicated types of models, we formulated a generalized Bayesian strategy for

₁₈₃ parameter estimation (see Appendix S1). For simplicity, we generated posterior predictions
₁₈₄ of ribbon seal abundance across the landscape as

$$N_i \sim \text{Poisson}(A_i \lambda_i), \quad (10)$$

₁₈₅ where A_i gives the proportion of suitable habitat in cell i (ribbon seals do not inhabit land
₁₈₆ masses).

₁₈₇ Fitting this model to our data,

₁₈₈ *Steller sea lion trends*

₁₈₉ DISCUSSION

₁₉₀ A number of authors have explored optimal knot placement in spatial models. In the
₁₉₁ context of predictive process modeling (where a covariance function is specified over a
₁₉₂ group of knots; see Banerjee et al., 2008), Finley et al. (2009) and Gelfand et al. (2013)
₁₉₃ considered near-optimal knot placement by minimizing spatially averaged prediction
₁₉₄ variance. Gelfand et al. knot selection

₁₉₅ Using GAs to select knot placement using

₁₉₆ Contrast with posterior loss (e.g., Jay's linex loss function estimator)

₁₉₇ Contrast with cross validation

₁₉₈ David Miller (the evil one)/Simon Wood stuff on edge effects

₁₉₉ Contrast with other approaches- Gelfand et al. Bayesian analysis - intrinsic CAR in
₂₀₀ SDM Chakraborty et al. '10 - spatial abundance modeling - ordinal Latimer et al. 2009 -
₂₀₁ spatial predictive process modelling in SPDs

202 Can't be complacent... still possible to get poor/biased results, e.g. if $\tau_\epsilon \rightarrow 0$. Can't
203 resolve pathological problems.

204 Presence-absence data, other link functions (e.g. probit)
205 extensions to models w/ secondary observation process, measurement error
206 much attention has been given to collinearity in multiple linear regression - suggest
207 researchers give as much attention to predictive extrapolation bias in predictive models
208 For predictions with spatial , our experience is that predictions outside the minimum
209 convex polygon where data are obtained can sometimes be more problematic than
210 predictions within the polygon. Spatial prediction surfaces may have a tendency to bend
211 up or down in these areas, resulting in “edge effects” that can lead to positive prediction
212 bias when a log link function is employed (Ver Hoef and Jansen, 2007).

213 Fieberg PVA

214 LITERATURE CITED

- 215 Banerjee, S., A. E. Gelfand, A. O. Finley, and H. Sang. 2008. Stationary process
216 approximation for the analysis of large spatial datasets. Journal of the Royal Statistical
217 Society B **70**:825–848.
- 218 Conn, P. B., J. M. Ver Hoef, B. T. McClintock, E. E. Moreland, J. M. London, M. F.
219 Cameron, S. P. Dahle, and P. L. Boveng. Accepted. Estimating multi-species abundance
220 using automated detection systems: ice-associated seals in the eastern Bering Sea.
221 Methods in Ecology and Evolution .
- 222 Cook, R. D. 1979. Influential observations in linear regression. Journal of the American
223 Statistical Association **74**:169–174.

- ²²⁴ Draper, N. R., and H. Smith. 1966. Applied Regression Analysis. John Wiley & Sons, New
²²⁵ York.
- ²²⁶ Ehrenberg, A., and J. Bound. 1993. Predicability and prediction. *Journal of the Royal*
²²⁷ Statistical Society A **156**:167–206.
- ²²⁸ Finley, A. O., H. Sang, S. Banerjee, and A. E. Gelfand. 2009. Improving the performance
²²⁹ of predictive process modeling for large datasets. *Computational Statistics & Data*
²³⁰ *Analysis* **53**:2873–2884.
- ²³¹ Gelfand, A. E., S. Banerjee, and A. O. Finley, 2013. Spatial design for knot selection in
²³² knot-based dimension reduction models. *in* J. Mateu and W. G. Muller, editors.
²³³ Spatio-temporal Design: Advances in Efficient Data Acquisition. Wiley.
- ²³⁴ Hastie, T. J., and R. J. Tibshirani. 1999. Generalized Additive Models. Chapman &
²³⁵ Hall/CRC, Boca Raton, Florida.
- ²³⁶ Lichstein, J., T. Simons, S. Shriner, and K. E. Franzreb. 2002. Spatial autocorrelation and
²³⁷ autoregressive models in ecology. *Ecological Monographs* **72**:445–463.
- ²³⁸ McCullagh, P., and J. A. Nelder. 1989. Generalized Linear Models. Chapman and Hall,
²³⁹ New York.
- ²⁴⁰ Ver Hoef, J., and J. Jansen. 2007. Space-time zero-inflated count models of harbor seals.
²⁴¹ *Environmetrics* **18**:697–712.
- ²⁴² Ver Hoef, J. M., M. F. Cameron, P. L. Boveng, J. M. London, and E. E. Moreland. 2013. A
²⁴³ hierarchical model for abundance of three ice-associated seal species in the eastern Bering
²⁴⁴ Sea. *Statistical Methodology* page <http://dx.doi.org/10.1016/j.stamet.2013.03.001>.

²⁴⁵ Wood, S. N. 2006. Generalized additive models. Chapman & Hall/CRC, Boca Raton,
²⁴⁶ Florida.

²⁴⁷ FIGURE 1. A ribbon seal, *Phoca fasciata*; the focus of spatial modeling efforts in this
²⁴⁸ paper.

²⁴⁹ FIGURE 2. Examples IVHs constructed from simulated data. In panels A and B, the
²⁵⁰ investigator plans to model a hypothetical (unmeasured) response variable using a linear
²⁵¹ regression model as a function of a single covariate, x , obtained at a number of design
²⁵² points (denoted with an “x”). Using x as a simple linear effect (A), only predictions less
²⁵³ than the minimum observed value of x or greater than the maximum value of x are outside
²⁵⁴ the IVH (shaded area), as scaled prediction variance in these areas (solid line) are greater
²⁵⁵ than the maximum scaled prediction variance for observed data (dashed line). Using both
²⁵⁶ linear and quadratic effects of x , some intermediate points are also outside the IVH;
²⁵⁷ predictions at these points should also be viewed with caution. Panels C & D show a more
²⁵⁸ complicated IVH when the investigator wishes to relate an unmeasured response variable
²⁵⁹ to linear and quadratic effects of two covariates, x and y , either without interactions (C) or
²⁶⁰ with interactions (D). Any potential predictions in the shaded area are outside of the IVH
²⁶¹ and should be viewed with caution.

²⁶² FIGURE 3. Aerial surveys over the Bering Sea in spring of 2012 (blue lines) overlayed
²⁶³ on a tessellated surface composed of 25km by 25km grid cells. Gray indicates land, and
²⁶⁴ colored pixels indicate ribbon seal encounters (yellow: 1-2 seals; orange: 3-4 seals; magenta:
²⁶⁵ 5-9 seals; red: 10-15 seals). On average, photographs covered approximately 2.6km^2 (0.4%)
²⁶⁶ of each surveyed grid cell.

²⁶⁷ FIGURE 4. Potential covariates gathered to help explain and predict ribbon seal
²⁶⁸ abundance in the eastern Bering Sea. Covariates include distance from mainland

269 (`dist_mainland`), distance from 1000m depth contour (`dist_shelf`), average remotely
270 sensed sea ice concentration while surveys were being conducted (`ice_conc`), and distance
271 from the southern sea ice edge (`dist_edge`). All covariates except ice concentration were
272 standardized to have a mean of 1.0 prior to plotting and analysis.

273 FIGURE 5 Posterior median estimates of ribbon seal apparent abundance across the
274 eastern Bering sea for (A) a generalized linear model (GLM), (B) a generalized additive
275 model (GAM), (C) a GLM with known zero data, and (D) a GAM with known zero data.
276 Highlighted cells indicate those where predictive covariate values are outside of the
277 generalized independent variable hull.

FIGURES



FIG 1

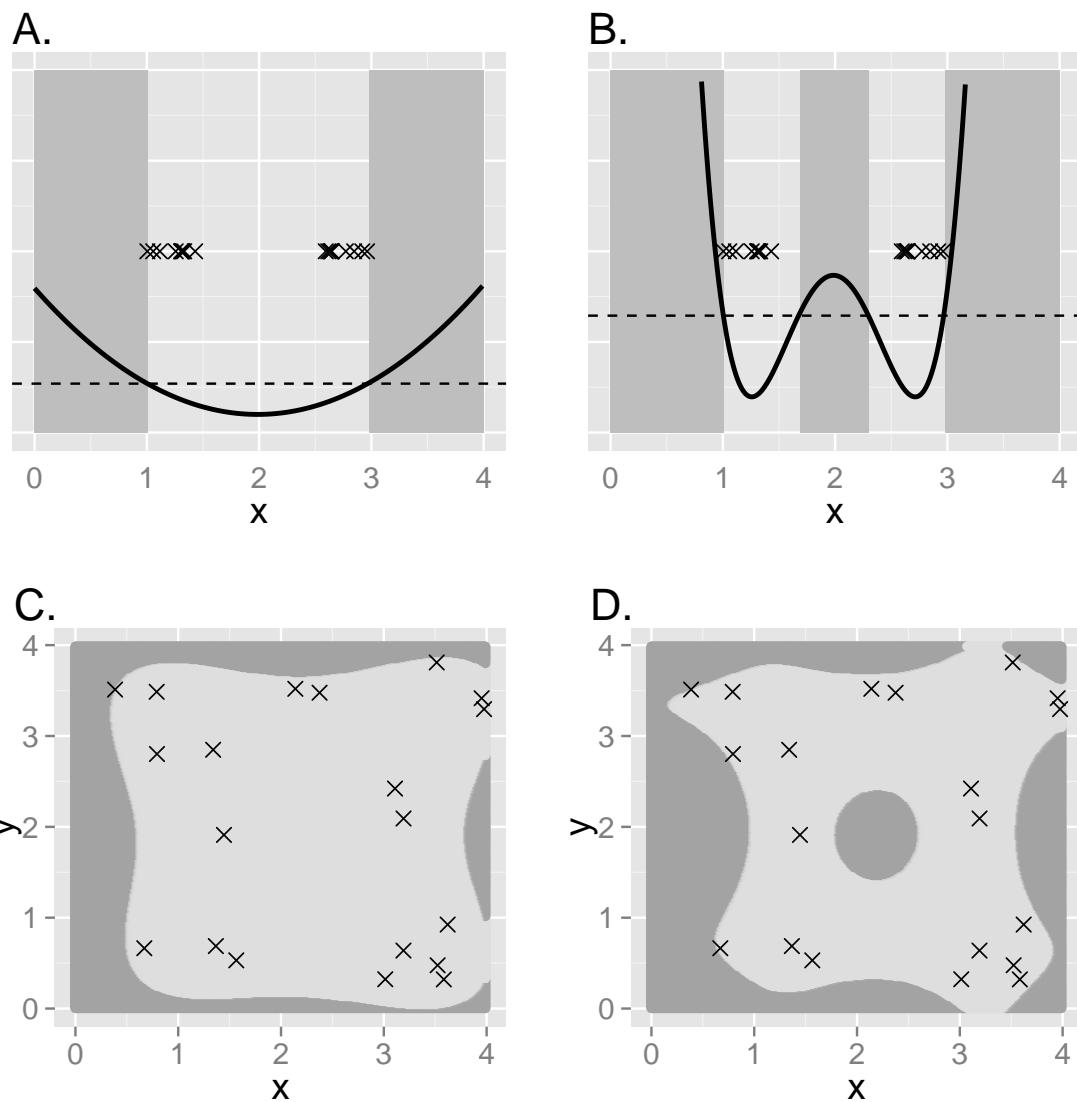


FIG 2

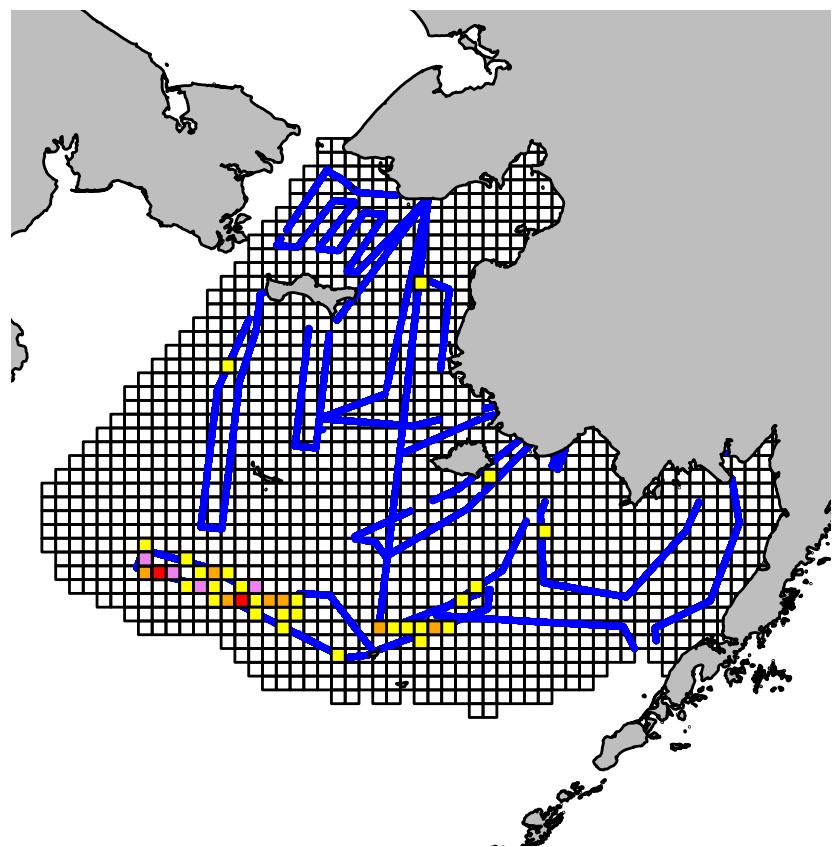


FIG 3

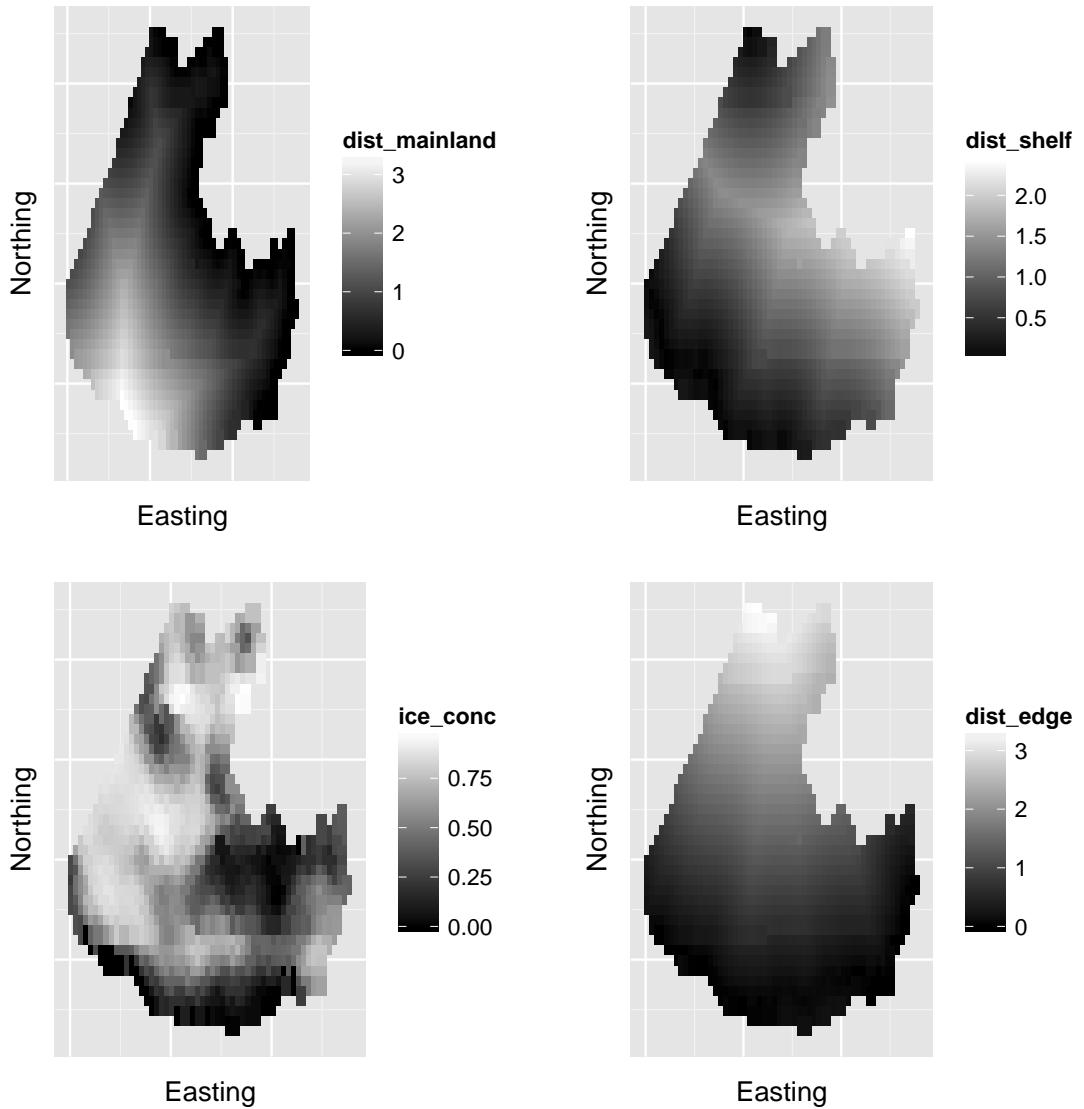


FIG 4

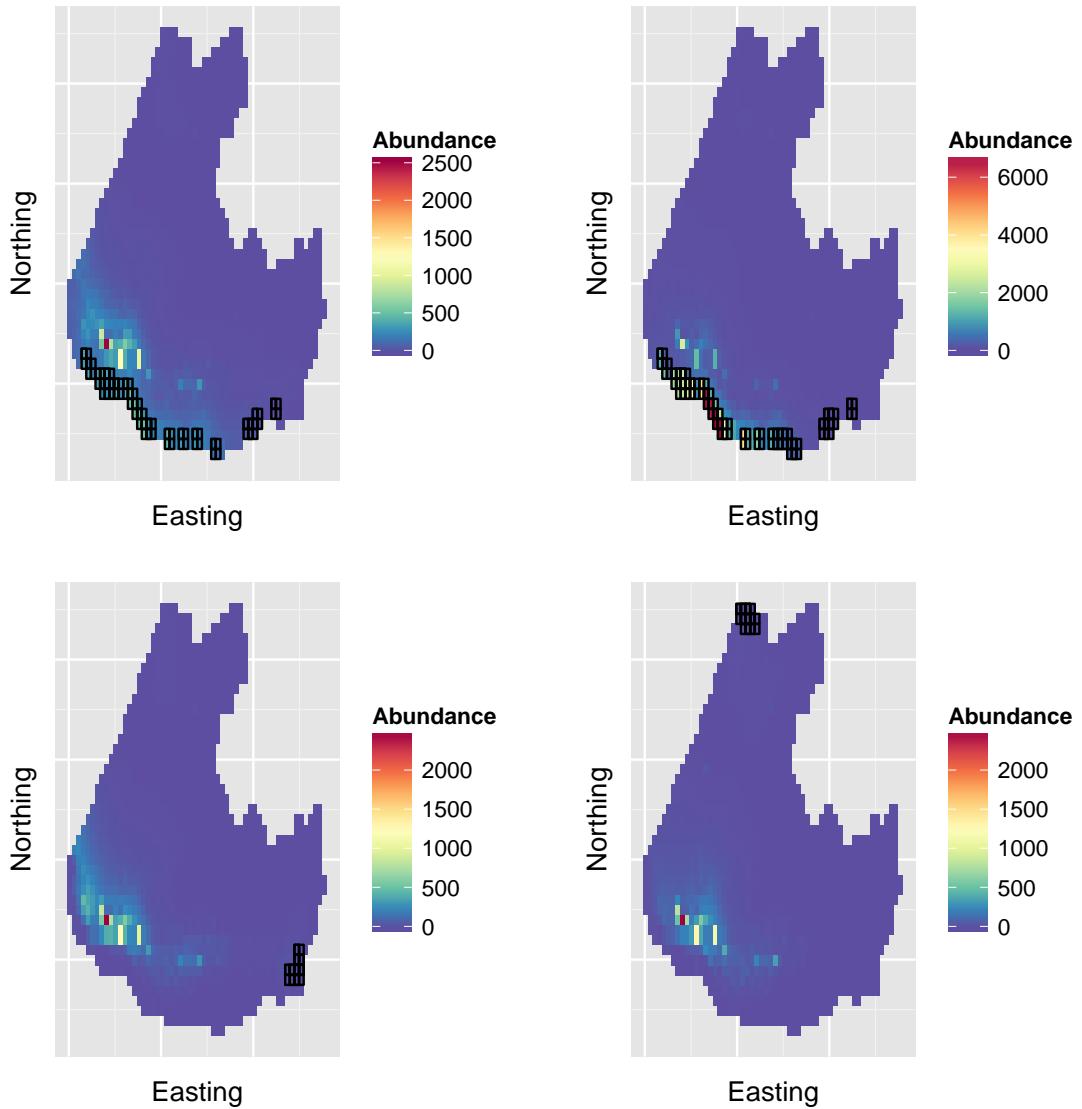


FIG 5