

RESEARCH ARTICLE

Combining multiple data sources in species distribution models while accounting for spatial dependence and overfitting with combined penalized likelihood maximization

Ian W. Renner¹  | Julie Louvrier²  | Olivier Gimenez³ 

¹School of Mathematical and Physical Sciences, The University of Newcastle, Callaghan, Australia

²Department of Ecological Dynamics, Department of Evolutionary Ecology, Leibniz Institute for Zoo and Wildlife Research, Berlin, Germany

³CEFE, CNRS, Univ Montpellier, Univ Paul Valéry Montpellier 3, EPHE, IRD, Montpellier, France

Correspondence

Ian W. Renner

Email: ian.renner@newcastle.edu.au

Funding information

University of Newcastle, Grant/Award Number: Early-to-Mid Career Researcher Visiting Fellowship; Université Paul-Valéry-Montpellier, Grant/Award Number: Professeurs en mobilité universitaire fund

Handling Editor: Robert B. O'Hara

Abstract

1. The increase in availability of species datasets means that approaches to species distribution modelling that incorporate multiple datasets are in greater demand. Recent methodological developments in this area have led to combined likelihood approaches, in which a log-likelihood comprised of the sum of the log-likelihood components of each data source is maximized. Often, these approaches make use of at least one presence-only dataset and use the log-likelihood of an inhomogeneous Poisson point process model in the combined likelihood construction. While these advancements have been shown to improve predictive performance, they do not currently address challenges in presence-only modelling such as checking and correcting for violations of the independence assumption of a Poisson point process model or more general challenges in species distribution modelling such as overfitting.
2. In this paper, we present an extension of the combined likelihood framework which accommodates alternative presence-only likelihoods in the presence of spatial dependence as well as lasso-type penalties to account for potential overfitting. We compare the proposed combined penalized likelihood approach to the standard combined likelihood approach via simulation and apply the method to modelling the distribution of the Eurasian lynx in the Jura Mountains in eastern France.
3. The simulations show that the proposed combined penalized likelihood approach has better predictive performance than the standard approach when spatial dependence is present in the data. The lynx analysis shows that the predicted maps vary significantly between the model fitted with the proposed combined penalized approach accounting for spatial dependence and the model fitted with the standard combined likelihood.
4. This work highlights the benefits of careful consideration of the presence-only components of the combined likelihood formulation, and allows greater flexibility and ability to accommodate real datasets.

KEYWORDS

area-interaction models, combined likelihood framework, diagnostic tools, lasso-type penalties, occupancy models, point process models, presence-only data, species distribution models

1 | INTRODUCTION

Species distribution models (SDMs), in which the distributions of species are modelled as a function of environmental predictors, rely on information about where a species has been observed (Guisan, Thuiller, & Zimmermann, 2017). Different SDM methods have been developed over the past few decades to accommodate the different protocols by which this species information is collected. For example, logistic regression and its extensions are often used when species detections and non-detections are recorded at a set of systematically designed locations (known as 'presence-absence' data), while point process models (PPMs, see Renner et al., 2015 for an overview) have emerged as a unifying framework for fitting SDMs informed by 'presence-only' data, in which only information about species presence locations is available. Statistically, these methods are often fitted by maximizing a corresponding likelihood expression, and the parameter estimates which maximize the likelihood may be used to produce maps of relative habitat suitability, reported as a habitat suitability index (Hirzel, Hausser, Chessel, & Perrin, 2002), probability of species presence (Phillips, Anderson, & Schapire, 2006), or intensity of locations per unit area (Warton & Shepherd, 2010) depending on the method.

Increasingly, species data are available from multiple sources and types. Many papers have advocated for fitting models to a combination of the available data types, illustrating benefits in model performance (Miller, Pacifici, Sanderlin, & Reich, 2019). Dorazio (2014) illustrated via simulations that adding a small amount of systematically collected presence-absence data to available presence-only data significantly improves predictive performance. Fithian, Elith, Hastie, and Keith (2015) showed that fitting a combined presence-only and presence-absence model to multiple species leverages the information of more abundant species to improve predictive performance for less prevalent species and allows sampling bias inherent in presence-only data to be estimated and corrected. These models are fitted by maximizing a combined log-likelihood expression which is the sum of the log-likelihoods of the presence-only and presence-absence components:

$$\ell(\alpha, \beta; \mathbf{s}_{PO}, \mathbf{y}_{PA}) = \ell_{PO}(\alpha_{PO}, \beta; \mathbf{s}_{PO}) + \ell_{PA}(\alpha_{PA}, \beta; \mathbf{y}_{PA}).$$

Here, \mathbf{s}_{PO} contains the locations of a presence-only data source, while \mathbf{y}_{PA} contains a vector of presence-absence detections and non-detections at a set of pre-selected sites. Parameters associated with the observation process unique to the presence-only and presence-absence datasets are denoted by α_{PO} and α_{PA} , respectively, and collectively contained in the vector α . Hereafter, we refer to these parameters as sampling bias parameters, as they may bias the intensity estimates as a result of the process of sampling the data. The key advancement of the combined likelihood approach is that the environmental response, parameterized by β , is informed by both the presence-only and presence-absence data.

Such an approach implicitly assumes that the datasets are statistically independent, which allows for the combined log-likelihood to be expressed as a sum of the single-source log-likelihoods.

Other combinations may be done in similar fashion. For example, Koshkina et al. (2017) considered a combination of presence-only data with site-occupancy data, and Pacifici et al. (2017) developed a multivariate conditional autoregressive model to account for spatial autocorrelation in occurrence and detection error.

While these papers clearly advance the practice of fitting SDMs in important ways, they do not address some common challenges that arise in real datasets. For example, they all consider an inhomogeneous Poisson point process model (IPPPM) for the presence-only data in the combination. In many real datasets, however, the implicit assumption that the point locations are independently distributed conditional on the environment is not met. Residual clustering or repulsion of the point locations not accounted for with an IPPPM due to the observation process, unconsidered environmental covariates, or biological factors would hence render the IPPPM inappropriate. One option to account for spatial dependence is to consider a log-Gaussian Cox Process, as Gelfand and Shirota (2018) do for a combination of presence-only and presence-absence data. Furthermore, none of the current literature in combined likelihood approaches includes ways to account for possible overfitting that results from including too many covariates in the model.

However, advances in SDM literature provide solutions to these common problems. Diagnostic tools such as the inhomogeneous K -function (Baddeley & Turner, 2000) and its simulation envelope (Diggle, 2003) can be used to determine departures from the independence assumption, and a wide number of alternative PPMs which account for spatial dependence may be included in the likelihood combination instead. Furthermore, penalized regression techniques such as the lasso penalty (Tibshirani, 1996) and its extension the adaptive lasso (Zou, 2006) may be used as a way to perform variable selection. Lasso regularization has been shown to boost predictive performance of SDMs and has been applied to IPPPMs (Renner & Warton, 2013) and occupancy models (Hutchinson, Valente, Emerson, Betts, & Dietterich, 2015).

In this paper, we present a penalized combined likelihood model in a way that it is more suitable for real datasets. In particular, we accommodate alternative forms of presence-only models to account for spatial dependence and affix a penalty on model complexity to address overfitting. In Section 2, we present the penalized combined likelihood formulation. In Section 3, we illustrate via simulations the improvements that this formulation provides and apply the proposed formulation to analyse the distribution of the Eurasian lynx *Lynx lynx* in the Jura Mountains in eastern France. Finally, we present a discussion and further avenues for research in this area in Section 4.

2 | MATERIALS AND METHODS

2.1 | Combined penalized likelihood formulation

We define the weighted, combined penalized log-likelihood as follows:

$$\ell(\alpha, \beta; \mathbf{y}) = \sum_{i=1}^D \ell_i(\alpha_i, \beta; \mathbf{y}_i) - p(\alpha, \beta). \quad (1)$$

Here, $\alpha = (\alpha_1, \dots, \alpha_D)^T$ is a q -dimensional vector that collects coefficients for the variables \mathbf{Z} used to model sampling bias for each of the D components individually. The environmental response is measured by a set of variables \mathbf{X} and is parametrized by $\beta = (\beta_1, \dots, \beta_p)^T$, which is collectively informed by all D components. The species data for all D components is collected in a set \mathbf{y} , with each individual data source \mathbf{y}_i determining the form of the component likelihood $\ell_i(\alpha_i, \beta; \mathbf{y}_i)$. Finally, $p(\alpha, \beta)$ is a penalty term described in further detail below.

While many possibilities for the likelihood terms $\ell_i(\alpha_i, \beta; \mathbf{y}_i)$ are possible, we will focus on likelihood expressions for a PPM and for an occupancy model. For an IPPPM, we typically model the intensity of points $\mu(s)$ over a given study region \mathcal{A} as a log-linear function of environmental variables \mathbf{X} and sampling bias terms \mathbf{Z} and derive estimates $\hat{\beta}$ and $\hat{\alpha}_{PO}$ of the associated parameters by maximizing a log-likelihood expression given by (Cressie, 1992):

$$\ell_{PO}(\alpha_{PO}, \beta; \mathbf{s}_{PO}) = \sum_{s \in \mathbf{s}_{PO}} \ln \mu(s) - \int_{s \in \mathcal{A}} \mu(s) ds. \quad (2)$$

In the simple occupancy model we consider, each site i is visited J_i times. We collect the history of detections and non-detections for all N sites in a matrix \mathbf{y}_{occ} . We assume that the probability that site i is occupied is given by ψ_i and that the occupancy of the sites remains constant throughout the history of visits. We further assume the probability of detecting the species if present is p_i . Under these assumptions, we can then model the probability of observing y_i detections at site i as

$$P(Y_i = y_i) = \underbrace{\psi_i \binom{J_i}{y_i} p_i^{y_i} (1 - p_i)^{J_i - y_i}}_{\text{species present}} + \underbrace{I(y_i = 0)(1 - \psi_i)}_{\text{species absent}},$$

where $I(\cdot)$ is the indicator function.

We can relate the occupancy ψ_i of site i to an inhomogeneous Poisson intensity μ_i of the species distribution over site i as in Koshkina et al. (2017):

$$\psi_i = 1 - e^{-\mu_i \times A_i},$$

where A_i is the area of site i . Note that $\mu_i \times A_i$ is an approximation of $\int_{s \in \text{site } i} \mu(s) ds$ that is reasonable if μ_i reasonably approximates the average intensity within site i .

As with the IPPPM, we can then model intensity as a log-linear function of environmental variables \mathbf{X} and model detection probability p_i as a function of some detection covariates \mathbf{Z} , such as the logit or complementary log-log function. We can then compute estimates $\hat{\beta}$ and $\hat{\alpha}_{occ}$ of the associated model parameters by maximizing the log-likelihood expression given by the following:

$$\ell_{occ}(\alpha_{occ}, \beta; \mathbf{y}_{occ}) = \ln \prod_{i=1}^N P(Y_i = y_i).$$

The term $p(\alpha, \beta)$ in Equation 1 is a penalty on model complexity applied to both the environmental parameters β and the sampling

bias parameters α to shrink these parameters toward zero in order to boost predictive performance. Here, we consider both the traditional lasso penalty (Tibshirani, 1996) and the adaptive lasso penalty (Zou, 2006). For the traditional lasso penalty,

$$p(\alpha, \beta) = \lambda \left(\sum_{j=1}^p |\beta_j| + \sum_{k=1}^q |\alpha_k| \right),$$

where λ is the tuning parameter. For the adaptive lasso penalty,

$$p(\alpha, \beta, \gamma) = \lambda \left(\sum_{j=1}^p w_j |\beta_j| + w_{p+k} \sum_{k=1}^q |\alpha_k| \right),$$

where $\mathbf{w} = (w_1, \dots, w_{p+q})^T$ are weights for the adaptive lasso, typically of the form:

$$w_i = \begin{cases} |\hat{\beta}_i^{(unp)}|^{-\gamma} & 1 \leq i \leq p \\ |\hat{\alpha}_{i-p}^{(unp)}|^{-\gamma} & p+1 \leq i \leq p+q, \end{cases}$$

for $\gamma > 0$. Here, $\hat{\beta}_i^{(unp)}$ is the unpenalized coefficient estimate corresponding to the i th environmental variable \mathbf{x}_i and $\hat{\alpha}_{i-p}^{(unp)}$ is the unpenalized coefficient estimate corresponding to the i th sampling bias variable \mathbf{z}_i . The shape of the weights is determined by the parameter γ . The data-driven choice of the adaptive weights \mathbf{w} ensures that more important covariates (i.e. those with coefficient estimates further away from 0) will be penalized less. This construction also enables the adaptive lasso to achieve the so-called oracle properties (Zou, 2006), which means that asymptotically, the correct subset of coefficients will be chosen and the procedure has an optimal estimation rate.

We can use Equation 1 to represent the simpler framework introduced by Dorazio (2014) and Fithian et al. (2015) by setting $p(\alpha, \beta) = 0$. We further extend this framework by considering alternative choices for those component likelihoods $\ell_i(\alpha_i, \beta; \mathbf{y}_i)$ informed by presence-only data. Rather than consider only inhomogeneous Poisson point process models, we consider area-interaction models (Baddeley & van Lieshout, 1995; Widom & Rowlinson, 1970) when diagnostic analysis of these data sources identifies spatial dependence among the presence-only locations. Area-interaction models account for spatial dependence through a vector of computed point interactions \mathbf{t}_s , which measure the proportion of overlap among circles of a nominal radius around the observed points \mathbf{s} . They can account for both clustering and repulsion of points – the model parameter η characterizes the nature of the spatial dependence, with values of η less than 1 signalling point repulsion and values of η greater than 1 signalling point clustering.

Because the likelihood expression of an area-interaction model is intractable, it is typically fitted via maximum pseudolikelihood (Besag, 1977):

$$\ell_{AI}(\alpha_{PO}, \beta, \eta; \mathbf{s}_{PO}) = \sum_{s \in \mathbf{s}_{PO}} \ln \mu(s; \mathbf{s}_{PO}) - \int_{s \in \mathcal{A}} \mu(s; \mathbf{s}_{PO}) ds.$$

This log-pseudolikelihood expression appears the same as Equation 2, with the exception that the intensity $\mu(s)$ is replaced by conditional intensity $\mu(s; s_{PO})$ (Papangelou, 1974), reflecting the fact that for the area-interaction model, intensity at a location s is conditional on the other points in the pattern s_{PO} .

2.2 | Implementation in R

To fit models with the combined penalized log-likelihood in Equation 1, we have developed a set of functions in R inspired by the `optim` function and `ppmlasso` package (Renner & Warton, 2013). The main function `comb_lasso` takes as an input a list of species data, associated environmental data, and formulae for the environmental trend and sampling bias trends for each component, along with details such as type of presence-only likelihoods to use the type of penalty, the number of models to fit, and the tuning parameter criterion. The function applies the coordinate descent algorithm of Osborne, Presnell, and Turlach (2000). This requires the derivatives of the component likelihoods (also known as 'score equations') to be computed. Analytical score equations are supplied directly to the `optim` function, which serves as the machinery of the optimization. A tutorial illustrating use of this code for the simulations as performed in Section 3.1 as well as some functions written to plot intensity maps and features of the lasso penalization is provided in Supporting Information.

3 | RESULTS

3.1 | Simulations

To investigate the benefits of the proposed penalized combined likelihood formulation, we used the `rpoispp` function in `spatstat` (Baddeley & Turner, 2005) to generate a large inhomogeneous Poisson pattern s_{true} of roughly 10,000 points on a 30×30 -unit square window from an intensity pattern defined by linear and quadratic terms of two generated variables (hence four meaningful covariates x_1, \dots, x_4 parameterized by coefficients β_1, \dots, β_4).

From this pattern, we generated two presence-only subsamples s_1 and s_2 biased by a different observation process. The first presence-only subsample s_1 was biased by z_1 , the distance to a simulated road network, and the other s_2 by z_2 , the distance to a simulated categorical covariate. We varied the size of the subsamples such that each pattern had 25, 100, or 400 points. We also varied the strength of the clustering of the presence-only subsamples by setting the coefficient of the interaction term $v_i = \ln \eta_i$ for $i = 1, 2$. Here, the patterns either exhibit no clustering ($v_i = 0$), moderate clustering ($v_i = 0.5$) or strong clustering ($v_i = 1$). In each case, the radius of interactions is set to 1 spatial unit. To sample the points in s_1 , we proceed as follows:

1. Initialize the set of sampled points $s_1 = \emptyset$ and the point interactions t_{s_1} to be a vector of 0s.
2. Compute the biased conditional intensity $\mu_1(s; s_1)$ at every point in s_{true} using x_1, \dots, x_4 , the sampling bias covariate z_1 , and the current

vector of point interactions t_{s_1} , where the biased conditional intensity is defined as follows:

$$\ln \mu_1(s; s_1) = \beta_1 x_1(s) + \beta_2 x_2(s) + \beta_3 x_3(s) + \beta_4 x_4(s) + \alpha_1 z_1(s) + v_1 t_{s_1}(s)$$

3. Set $\mu_1(s; s_1) = 0$ for all $s \in s_1$. That is, we set the conditional intensity for any point already selected in s_1 to 0 to ensure these points are not resampled.
4. Randomly select a point from s_{true} with sampling probabilities proportional to the conditional intensities and add the selected point to s_1 .
5. Update the vector of point interactions t_{s_1} for all points in s_{true} using the internal `evalInteraction` function in `spatstat`, which computes point interactions based on a supplied point pattern for a given set of locations and interaction radius.
6. Repeat steps 2–5 until we have sampled the desired number of points.

We sample s_2 in a similar manner, using z_2 instead of z_1 to create the sampling bias and computing point interactions t_{s_2} .

Because the true pattern s_{true} is Poisson, this simulation setup emulates a scenario in which the clustering of the observed point patterns is an artefact of the observation process – this can happen if, for example, records are publicly available and enthusiasts for the species report further observations near the publicly available locations (Johnston et al., 2019).

We also generated a history y_{occ} of detections and non-detections from 5 visits to each of 100 sites centred along a regular grid in the 30×30 -unit observation window to emulate a dataset for which we could consider occupancy modelling. The species was considered present at a site if the closest point in the pattern s_{true} was within a distance of 0.18 units of the centre of the site, such that the area of each site is roughly 0.1 square units. The history of detections and non-detections at each site where the species was considered present was randomly generated according to detection probabilities defined by the inverse of the cloglog function evaluated at a generated detection covariate z_3 .

Finally, we generated four dummy covariates d_1, \dots, d_4 to include in fitted models that were meaningless in describing the true species distribution. We did this to reflect the fact that in real applications, we may not know which among a suite of candidate variables truly determine the species distribution. We ensured that the maximum absolute correlation among all pairs of variables was smaller than 0.5.

After generating the species data, we fit a number of models, using as input environmental covariates the four meaningful covariates x_1, \dots, x_4 (parameterized by β_1, \dots, β_4) as well as four dummy covariates d_1, \dots, d_4 (parameterized by β_5, \dots, β_8) and using as sampling bias covariates z_1, z_2 , and z_3 (parameterized by α_1, α_2 , and α_3). For both Poisson and area-interaction presence-only likelihoods, we fit a model without any penalty, with a lasso penalty, and with an

adaptive lasso penalty. For the models fitted with either a lasso or an adaptive lasso penalty, we fit regularization paths of 1,000 models, increasing the penalty from 0 to the smallest penalty λ_{\max} that would shrink all coefficients to 0, thus covering the entire scope of possible model sizes. The model which minimized BIC was chosen among the 1,000 fitted models. We considered as species data a combination of all three of s_1 , s_2 , and y_{occ} . This led to a total of six models being fitted, summarized in Table 1.

To evaluate performance, we compared the integrated mean squared error of the true intensity surface with rescaled fitted intensity surfaces of the six models. The fitted intensity surfaces were rescaled to have the same mean intensity as the true intensity surface to ensure that fair comparisons are made as models using different species data sources will have varying intercepts to reflect the estimated abundance of the points.

We performed 1,000 simulations of the datasets for each of the nine combinations of presence-only dataset size and clustering strength and the resultant model fits on 512 GB nodes powered by 3.0 GHz Intel Xeon Gold (E5-6154) processor from the University of Newcastle's High Performance Computing cluster. The 9,000 simulation tasks took approximately 7,000 hr.

Figure 1 shows boxplots of the calculated integrated mean squared errors from the simulations. From these results, we can draw the following conclusions. First, the models fitted with the area-interaction presence-only likelihoods have performance benefits over the models fitted with Poisson presence-only likelihoods when clustering is present. When clustering is not present (first column), a setting for which the Poisson likelihood is appropriate, the models fitted with area-interaction presence-only likelihoods perform no worse than models fitted with Poisson presence-only likelihoods. Comparing the plots across rows and down columns, we see that the performance advantage of the models fitted with area-interaction presence-only likelihoods tends to increase as the degree of clustering gets larger and as the sample size increases, respectively.

In Appendix S1, we show that the parameter coefficients β_1, \dots, β_4 corresponding to the meaningful covariates x_1, \dots, x_4 are increasingly biased away from 0 for the models fitted with Poisson presence-only likelihoods, both as sample size increases and as the strength of presence-only clustering increases. The inclusion of the area-interaction

term takes an increasing amount of signal from the environmental covariates as the strength of the presence-only clustering increases. For low sample sizes, there is a suggestion that this signal dampening may be too strong, though such an overcorrection disappears as sample size increases.

Second, penalization via the lasso or adaptive lasso improves model performance when there is no presence-only clustering, and this improvement is greatest for smaller sample sizes. This is an expected conclusion given the danger of overfitting is greater with fewer observations. Models penalized with the adaptive lasso tend to outperform models penalized with the lasso when there is no presence-only clustering. However, lasso penalization does not notably improve performance when there is presence-only clustering. In fact, there is a suggestion that applying a lasso penalty may slightly hinder performance when applying an area-interaction presence-only likelihood for small sample sizes. Although the benefits of penalization are negligible with large datasets, fitting models with a penalty does not hurt the performance.

In summary, it appears that the proposed combined penalized likelihood framework provides the best performance. Furthermore, incorporating area-interaction presence-only likelihoods improves performance when clustering is present, and can likewise reliably estimate that there are negligible point interactions if clustering is not present, in effect relaxing to the simpler model with Poisson presence-only likelihoods when this additional complexity is not needed. A more detailed discussion of the simulation results, including boxplots of the fitted coefficients, appears in Appendix S1.

3.2 | Analysis of Eurasian lynx distribution in the Jura Mountains

We now demonstrate the use of the combined penalized likelihood approach to analyse the distribution of the Eurasian lynx in the Jura Mountains in eastern France.

Lynx went extinct in France at the end of the 19th century due to habitat degradation, human persecution and decrease in prey availability (Vandel & Stahl, 2005). The species was reintroduced in Switzerland in the 1970s (Breitenmoser, Breitenmoser-Würsten, & Capt, 1998), then re-colonized France through the Jura mountains in the 1980s (Vandel & Stahl, 2005). The species is listed as endangered under the 2017 IUCN Red list and is of conservation concern in France due to habitat fragmentation, poaching and collisions with vehicles. The Jura holds the bulk of the French lynx population.

We have three sources of lynx data in the Jura Mountains: a presence-only dataset consisting of 440 opportunistic sightings in the wild from 2009–2011 (denoted s_w), another presence-only dataset consisting of 240 reported interferences of lynx with domestic livestock in 2009–2011 (denoted s_d), and pictures of lynx taken from cameras set up in 73 locations s_c in the Jura Mountains in 2012. Lynx presence-only data were made of presence signs sampled all year long thanks to a network of professional and non-professional observers. Every observer is trained during a 3-day teaching course led by the French National Game and Wildlife Agency (ONCFS) to

TABLE 1 Models fitted in each simulation using the proposed combined penalized likelihood. The models also varied based on the likelihood expression for any presence-only components and the type of penalty used, if any

| Model | Species data | Presence-only likelihood | Penalty |
|-------|-----------------------------------|--------------------------|----------------|
| 1 | s_1, s_2 , and y_{occ} | IPPPM | None |
| 2 | s_1, s_2 , and y_{occ} | IPPPM | Lasso |
| 3 | s_1, s_2 , and y_{occ} | IPPPM | Adaptive Lasso |
| 4 | s_1, s_2 , and y_{occ} | Area-interaction | None |
| 5 | s_1, s_2 , and y_{occ} | Area-interaction | Lasso |
| 6 | s_1, s_2 , and y_{occ} | Area-interaction | Adaptive Lasso |

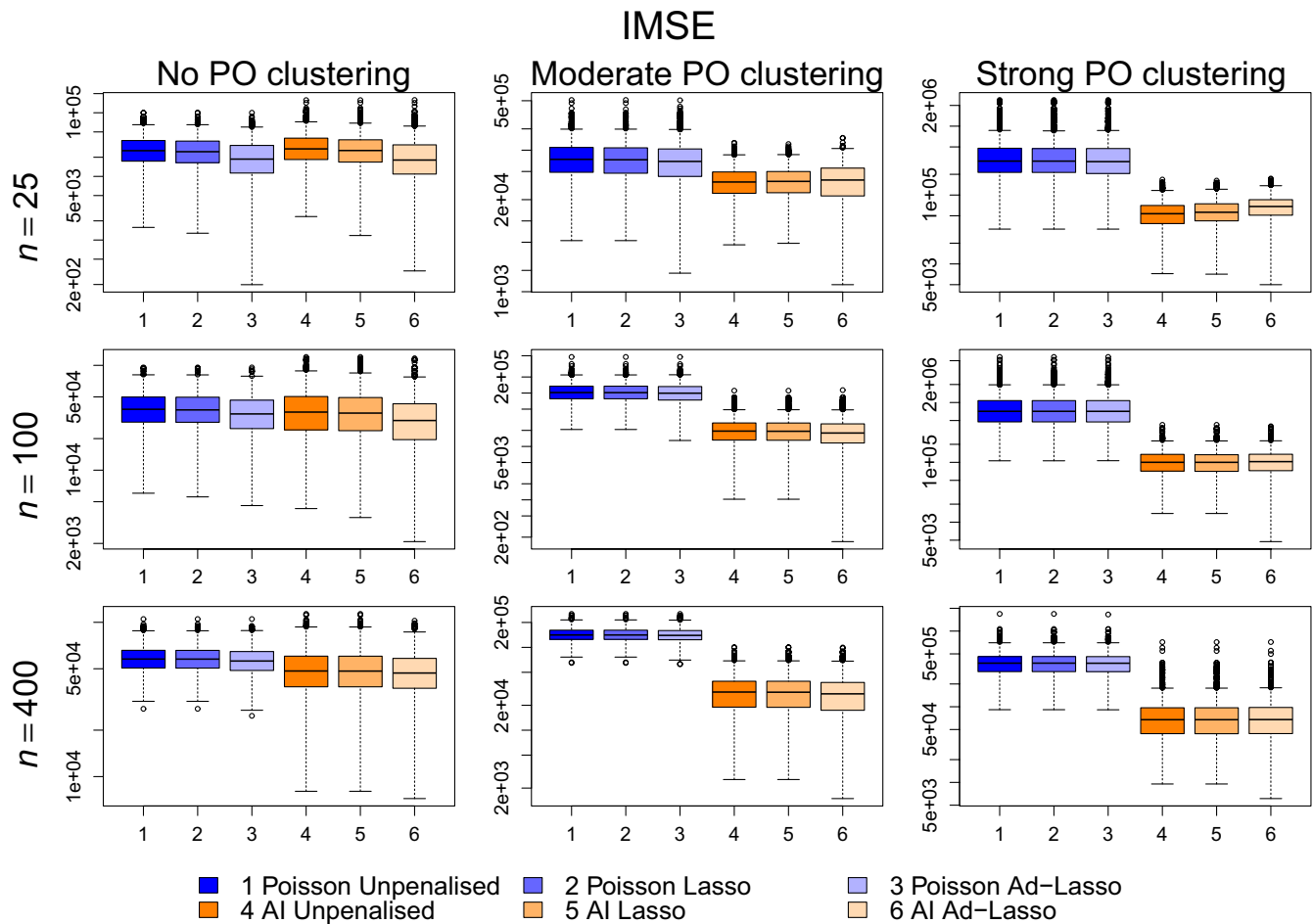


FIGURE 1 Boxplots of integrated mean squared error for the six models described in Table 1 for different combinations of presence-only sample size and clustering strength

document signs of the species' presence (Duchamp et al., 2012). Presence signs went through a standardized control process to prevent misidentification (Duchamp et al., 2012). The camera data have daily reportings of the lynx across a total of 77 days. Due to this, we can consider the picture history of lynx at the camera locations in an occupancy modelling framework (Blanc, Marboutin, Gatti, Zimmermann, & Gimenez, 2014). In particular, we split the 77-day period into seven 11-day periods, such that the site history y_c comprises seven detections and non-detections at each site in s_c over each 11-day period.

Figure 2 shows the locations of the sightings in both presence-only datasets as well as the locations of the cameras. Both presence-only data sources appear to have different distributions, reflecting different sampling biases. There are more wild sightings in the north-east of the Jura Mountains, and more domestic interferences toward the southwest. Additionally, there appear to be some tight clusters within both datasets, with several records very close to each other.

To model the lynx distribution, we consider altitude, percentage of forest cover, distance to the nearest water source and human population density as environmental variables. We model sampling bias in the wild records s_w with distance to the nearest main road and distance to the nearest train line, and sampling bias in the domestic

records s_d with distance to the nearest farm and percentage of agricultural land. Finally, we model detection probability for the camera data with distance to the nearest urban area. We established this set of potential candidate environmental and detection variables based on previously studied species habitat preferences and detectability (Bouyer et al., 2015). The Corine Land Cover land-use repository from 2012 (Büttner, Soukup, & Kosztra, 2014) supplies a map of land coverage including urban areas, water areas, forest areas, farm areas and agricultural areas that was used to generate the percentage of forest areas and agricultural areas over $1 \text{ km} \times 1 \text{ km}$ cells as well as distances to the nearest urban area, water source, and farm. Altitude was averaged over $1 \text{ km} \times 1 \text{ km}$ cells from data available in the `raster` package in R, while human population density was averaged over $1 \text{ km} \times 1 \text{ km}$ cells taken from version 4 of the Gridded Population of the World data repository (Center for International Earth Science Information Network (CIESIN) – Columbia University, 2016). Distances from the nearest main road and railway were computed from shapefiles from Version 151 of the ROUTE 500 database, accessible at <http://professionnels.ign.fr/route500>.

We fitted initial separate IPPPMs to the wild records s_w and the domestic records s_d using linear, quadratic, and interaction terms for the four environmental covariates, and linear terms for the sampling

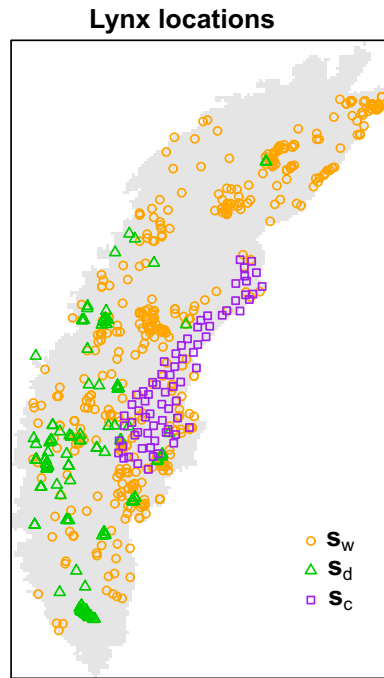


FIGURE 2 Locations of the lynx data in the Jura Mountains, including 440 observations in the wild s_w , 220 reports of domestic interference s_d , and 73 camera traps s_c

bias covariates. From these models, we are able to assess whether the assumption of independence inherent to the IPPPMs is appropriate with simulation envelopes of the inhomogeneous K -function in *spatstat*, as shown in Figure 3. Both of the envelopes for the IPPPMs fitted to the wild model (left panel) and the domestic model (middle panel) demonstrate additional clustering as the observed inhomogeneous K -function values plotted in red fall above the simulation envelopes for small radii. This suggests that fitting an IPPPM is inappropriate for these datasets. The right panel shows a simulation envelope of the cross K -function as produced by the `Kcross.inhom` function of *spatstat*, which counts the expected number of wild sightings within a given distance of a domestic sighting, conditional on the spatially varying intensities of both patterns. We estimate the wild and domestic intensities from area-interaction models, and as

the observed values of the cross K -function fall within the envelope boundaries, this suggests that there is no clustering across the two datasets. This, in turn, suggests that the observed clustering within the wild and domestic datasets may be more likely attributable to the observation process than to some biological reality that induces clustering or a missed environmental covariate.

Consequently, we fit combined likelihood models using both the standard, unpenalized approach (analogous to Model 1 in Table 1) and the combined penalized likelihood formulation Equation 1 with a lasso penalty and area-interaction models for the presence-only data sources (analogous to Model 5 in Table 1). The radii chosen to capture the residual spatial patterning in the wild and domestic models are 2 km and 5 km, as chosen by the `profilepl` function in *spatstat*.

Figure 4 shows the bias-corrected fitted intensities from these two models. For the combined model which uses IPPPMs (left panel), the fitted intensity is corrected for the sampling bias terms modelled for the presence-only components using the method of Warton, Renner, and Ramp (2013). For the combined penalized model which uses area-interaction models (right panel), the fitted intensity is corrected for these same sampling bias terms as well as the fitted point interactions – that is, we treat the interaction parameter ν as belonging to the set of sampling bias parameters α . The fitted models show strikingly different patterns, with the model which uses area-interaction components highlighting much more of the Jura Mountains as preferred habitat of lynx than the model which uses IPPPMs. The models suggest similar numbers of points throughout the Jura, but the distribution of these points are more heavily concentrated in the IPPPM model. This is because the area-interaction terms in the AI model lessen the impact of some clusters of points on the scale of the displayed bias-corrected intensities.

We do not have access to additional data with which to validate the performance of these models such as GPS data as in Gould, Gould, Cain, and Roemer (2019), but the results of Section 3.1 suggest that the model which uses area-interaction components is likely to better reflect the true distribution of lynx.

The combined penalized model with the area-interaction components found the optimal lasso penalty was 0, resulting in a model

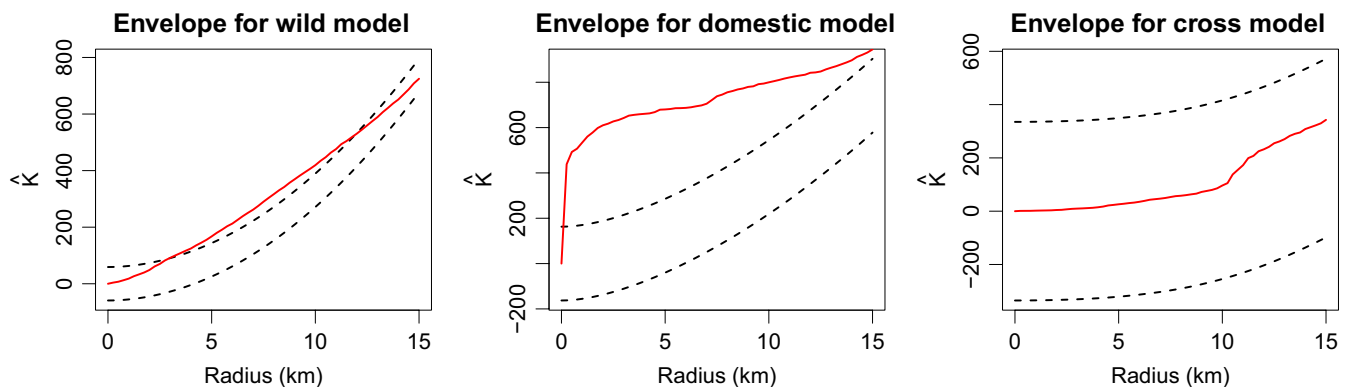
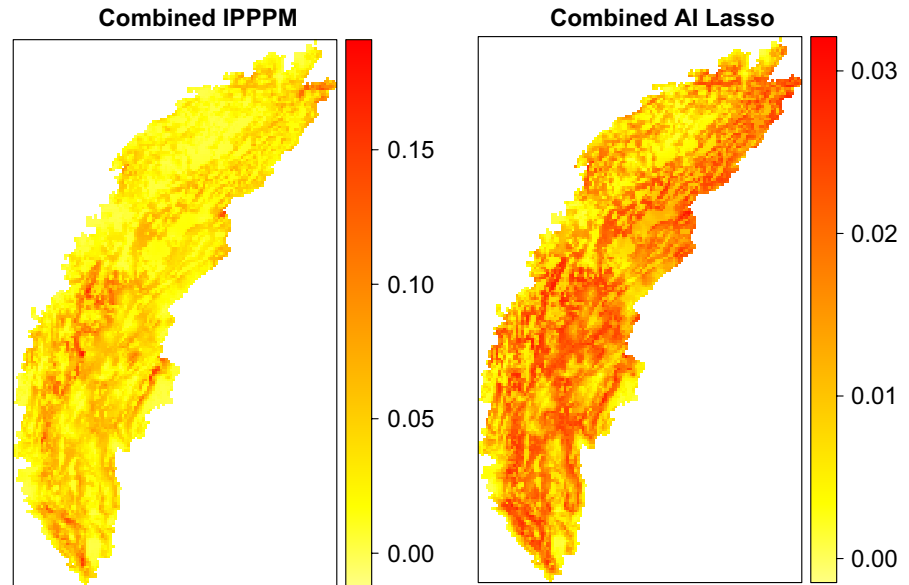


FIGURE 3 95% simulation envelopes of the inhomogeneous K -function for the fitted inhomogeneous Poisson point process model (IPPPM) of the wild records (left), the fitted IPPPM of the domestic records (middle), and across the two fitted IPPPMs (right)

FIGURE 4 Fitted intensities using the combined likelihood formulation. Left: the model is fitted without any penalty and using inhomogeneous Poisson point process models for the presence-only data sources. Right: the model is fitted with a lasso penalty and using area-interaction models for the presence-only sources



which included all 18 covariates and both of the area-interaction terms. The fact that the optimal penalty is 0 suggests that the suite of covariates we chose to include, motivated by existing literature, seems to have been a good choice. In general, we recommend use of the lasso penalty as a safeguard against overfitting, particularly in contexts where the suite of candidate covariates for a species is less established as an insurance against overfitting.

4 | DISCUSSION

The proposed combined penalized likelihood framework addresses some common problems that arise in real datasets. The flexibility to incorporate an area-interaction likelihood when there is spatial dependence in the presence-only dataset and affix a penalty on model complexity enables improvements in predictive performance, as shown in Section 3.1.

4.1 | Possible extensions

Despite these improvements, further advances are possible. Other penalty structures could be incorporated into the same framework. While the lasso and adaptive lasso showcased here show clear benefits in simulations, other penalized likelihood variants such as SCAD (Fan & Li, 2001) could lead to superior performance in some situations, and alternative methods to BIC of choosing the size of the penalty such as the Extended Bayesian Information Criterion ('ERIC', Hui, Warton, & Foster, 2015) could likewise be used.

While we make use of the area-interaction likelihood in this paper, there is a large family of Gibbs PPMs (Cressie, 1992) which accommodate different sorts of spatial dependence that could be used. Our choice of the area-interaction model as the alternative is motivated by the fact that it accommodates interactions of all orders instead of just pairwise interactions and that it can be used to model both clustering and repulsion of points.

The inclusion of the area-interaction terms dampens the signal of the environmental covariates. Although this makes sense when spatial dependence exists, we may dampen the signal too much. In the context of species distribution models, we might ask the question, 'Does a given species record exist because its location is in particularly suitable habitat for the species, or because there are other records nearby?' If the answer to this question appears to be 'both', as is often the case for presence-only data, we are at risk of 'spatial confounding'. In the single-source context, Hodges and Reich (2010) propose restricting the spatial effect to be orthogonal to the fixed covariate effects, while Simpson, Rue, Riebler, Martins, and Sørbye (2017) and Sørbye, Illian, Simpson, Burslem, and Rue (2019) suggest careful selection of associated spatial priors to alleviate this risk. With our implementation, we could achieve something similar to the latter two papers by adjusting the magnitude of the lasso penalty for the area-interaction terms. In Appendix S1, we highlight the tradeoff between the estimates of the interaction parameters $\hat{\nu}_i$ and both the estimates of the environmental parameters $\hat{\beta}_i$ and the sampling bias parameters $\hat{\alpha}_i$. However, a full exploration of the effects of spatial confounding remains an open area of research and is beyond the scope of this paper.

In both the simulations in Section 3.1 and the lynx data analysis in Section 3.2, we made the rather limiting assumption of a closed population and that sites are either always occupied or always unoccupied. Nonetheless, occupancy models which take into account changing site dynamics could be used (MacKenzie, Nichols, Hines, Knutson, & Franklin, 2003). Similarly, we have ignored the temporal aspect of the lynx distribution in this paper, but there is a wide suite of tools to fit spatio-temporal models in order to capture distribution dynamics for both the aforementioned occupancy modelling component as well as presence-only components (Cressie & Wile, 2015).

Further improvements could be made by incorporating source weights in situations in which the data sources vary in quality. Indeed, presence-only data sources may be more prone to errors in coordinate locations as well as correct species identification, as

they often include records by amateur enthusiasts. The combined penalized likelihood framework could easily be extended to include weights for the various data sources by adding a vector of source weights $\mathbf{w} = (w_1, \dots, w_D)^T$ to the formulation in Equation 1:

$$\ell(\alpha, \beta; \mathbf{y}) = \sum_{i=1}^D w_i \ell_i(\alpha_i, \beta; \mathbf{y}_i) - p(\alpha, \beta). \quad (3)$$

One possible strategy to incorporate such weights in Equation 3 could be to compare performance of single source models on independent data and upweight the contribution of data sources that are shown to have good performance.

Finally, while we incorporate sampling bias as a linear effect, nonlinear effects can also be used as appropriate for a given sampling protocol, for example with distance sampling as discussed in Yuan et al. (2017).

4.2 | Accounting for dependence within and among data sources

In the lynx data analysis in Section 3.2, we diagnosed spatial dependence within each of the presence-only data sources, but found no spatial dependence across data sources. Tools such as the inhomogeneous K -envelope provide great insight into the underlying individual spatial processes that are observed. However, such diagnostic tools are not currently available for the combined likelihood models, and research in this area would be valuable as these models grow in popularity.

Another approach to constructing SDMs from multiple data sources could be to introduce a common latent spatial term $\xi(\mathbf{s})$, such as a Gaussian random field, which would account for spatial dependence among points in all of the data sources as in Gelfand and Shirota (2018). The resulting likelihood expression would be:

$$\ell(\alpha, \beta; \mathbf{y}) = \sum_{i=1}^D \ell_i(\alpha_i, \beta; \mathbf{y}_i) + \xi(\mathbf{y}) - p(\alpha, \beta), \quad (4)$$

where $\xi(\mathbf{y}) \sim \text{MVN}(\mathbf{0}, \Sigma)$. Models of this type are typically fitted in a Bayesian framework. We could reduce the dimension of ξ through methods like fixed rank kriging or induce sparsity in Σ through lasso-type penalties such that the likelihood in Equation 4 could be fitted with software such as Template Model Builder (TMB, Kristensen, Nielsen, Berg, Skaug, & Bell, 2016). Another way to achieve sparsity is with the stochastic partial differential equation approach (SPDE, Lindgren, Rue, & Lindström, 2011), as implemented in the *inlabru* package (Bachl, Lindgren, Borchers, & Illian, 2019).

4.3 | Conclusion and perspectives

The development of statistical methods is often motivated by new challenges raised by novel types of datasets. While the current literature on combined likelihood approaches represents a significant recent advancement, advances in other areas can be lost

if not carried over with such methodological developments. This paper attempts to build a bridge between this exciting new arena for species distribution modelling and the rich suite of tools available for species distribution modelling, particularly that for presence-only data. Our hope is that other such bridges continue to be built in this spirit.

ACKNOWLEDGEMENTS

We thank the staff from the French National Game and Wildlife Agency, the Forest National Agency and the Departmental Federation of Hunters of Jura department, who collected the photographs during the camera-trapping session. We also thank all the volunteers who are members of the Réseau Loup-Lynx who collect every year precious presence signs of lynx all over its distribution area. We thank the University of Newcastle through the Early to Mid-Career Researcher Visiting Fellowship and the Université Paul-Valéry Montpellier through the Professeurs en mobilité universitaire fund for funding visits for I.W.R. which helped facilitate the collaborations that resulted in this work. O.G. was funded by CNRS and the 'Mission pour l'Interdisciplinarité' through the 'Osez l'Interdisciplinarité' initiative. Finally, we thank the reviewers for their very helpful comments, which have helped us greatly improve the paper.

AUTHORS' CONTRIBUTIONS

I.W.R. and O.G. conceived the concept of the paper. I.W.R. developed the code to fit the models. J.L. sourced the species coordinates and covariates for the lynx analysis. I.W.R. and O.G. wrote the manuscript. I.W.R. and J.L. developed the tutorial in Supporting Information. All authors were involved in editing drafts of the manuscript.

DATA AVAILABILITY STATEMENT

The Eurasian lynx is an endangered species with high conservation stakes. Interactions with human activities are problematic and lead to poaching and anthropogenic pressures. Providing accurate information on lynx locations can be detrimental to the conservation status of the species. As a consequence, the owners of the original data denied the request to archive it.

Interested readers may request access to the lynx data from Nolwenn Drouet-Hoguet, head of the lynx and unit at the Office National de la Chasse et de la Faune Sauvage, at nolwenn.drouet-hoguet@oncfs.gouv.fr.

ORCID

Ian W. Renner  <https://orcid.org/0000-0003-3116-2486>

Julie Louvrier  <https://orcid.org/0000-0003-1252-1746>

Olivier Gimenez  <https://orcid.org/0000-0001-7001-5142>

REFERENCES

- Bachl, F. E., Lindgren, F., Borchers, D. L., & Illian, J. B. (2019). INLABRU: An R package for Bayesian spatial modelling from ecological survey data. *Methods in Ecology and Evolution*, 10, 760–766. <https://doi.org/10.1111/2041-210X.13168>
- Baddeley, A., & Turner, R. (2000). Practical maximum pseudolikelihood for spatial point patterns (with discussion). *Australian & New Zealand Journal of Statistics*, 42, 283–322. <https://doi.org/10.1111/1467-842X.00128>
- Baddeley, A., & Turner, R. (2005). SPATSTAT: An R package for analyzing spatial point patterns. *Journal of Statistical Software*, 12, 1–42.
- Baddeley, A. J., & van Lieshout, M. N. M. (1995). Area-interaction point processes. *Annals of the Institute of Statistical Mathematics*, 47, 601–619. <https://doi.org/10.1007/BF01856536>
- Besag, J. (1977). Some methods of statistical analysis for spatial data. *Bulletin of the International Statistical Institute*, 47, 77–92.
- Blanc, L., Marboutin, E., Gatti, S., Zimmermann, F., & Gimenez, O. (2014). Improving abundance estimation by combining capture–recapture and occupancy data: Example with a large carnivore. *Journal of Applied Ecology*, 51, 1733–1739. <https://doi.org/10.1111/1365-2664.12319>
- Bouyer, Y., San Martin, G., Poncin, P., Beudels-Jamar, R. C., Odden, J., & Linnell, J. D. (2015). Eurasian lynx habitat selection in human-modified landscape in norway: Effects of different human habitat modifications and behavioral states. *Biological Conservation*, 191, 291–299. <https://doi.org/10.1016/j.biocon.2015.07.007>
- Breitenmoser, U., Breitenmoser-Würsten, C., & Capt, S. (1998). Re-introduction and present status of the lynx (*Lynx lynx*) in Switzerland. *Hystrix, The Italian Journal of Mammalogy*, 10, 17–30.
- Büttner, G., Soukup, T., & Kosztra, B. (2014). CLC2012 addendum to CLC2006 technical guidelines. Final Draft, Copenhagen (EEA).
- Center for International Earth Science Information Network (CIESIN) – Columbia University. (2016). Gridded population of the world, version 4 (gpwv4): Population density.
- Cressie, N. (1992). *Statistics for spatial data* (vol. 4). Wiley Online Library.
- Cressie, N., & Wile, C. K. (2015). *Statistics for spatio-temporal data*. New York, NY: John Wiley & Sons.
- Diggle, P. J. (2003). *Statistical analysis of spatial point patterns*. London: Edward Arnold.
- Dorazio, R. M. (2014). Accounting for imperfect detection and survey bias in statistical analysis of presence-only data. *Global Ecology and Biogeography*, 23, 1472–1484. <https://doi.org/10.1111/geb.12216>
- Duchamp, C., Boyer, J., Briaudet, P. E., Leonard, Y., Perrine Moris, P., Bataille, A. ... Marboutin, E. (2012). A dual frame survey to assess time-and space-related changes of the colonizing wolf population in France. *Hystrix-Italian Journal of Mammalogy*, 23, 14–28.
- Fan, J., & Li, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96, 1348–1360. <https://doi.org/10.1198/016214501753382273>
- Fithian, W., Elith, J., Hastie, T., & Keith, D. A. (2015). Bias correction in species distribution models: Pooling survey and collection data for multiple species. *Methods in Ecology and Evolution*, 6, 424–438. <https://doi.org/10.1111/2041-210X.12242>
- Gelfand, A., & Shirota, S. (2018) Preferential sampling for presence/absence data and for fusion of presence/absence data with presence-only data. *arXiv Preprint arXiv:180901322*.
- Gould, M. J., Gould, W. R., Cain, J. W. III, & Roemer, G. W. (2019). Validating the performance of occupancy models for estimating habitat use and predicting the distribution of highly-mobile species: A case study using the American black bear. *Biological Conservation*, 234, 28–36. <https://doi.org/10.1016/j.biocon.2019.03.010>
- Guisan, A., Thuiller, W., & Zimmermann, N. E. (2017). *Habitat suitability and distribution models: With applications in R*. Cambridge: Cambridge University Press.
- Hirzel, A. H., Hausser, J., Chessel, D., & Perrin, N. (2002). Ecological-niche factor analysis: How to compute habitat-suitability maps without absence data? *Ecology*, 83, 2027–2036. [https://doi.org/10.1890/0012-9658\(2002\)083\[2027:ENFAHT\]2.0.CO;2](https://doi.org/10.1890/0012-9658(2002)083[2027:ENFAHT]2.0.CO;2)
- Hodges, J. S., & Reich, B. J. (2010). Adding spatially-correlated errors can mess up the fixed effect you love. *The American Statistician*, 64, 325–334. <https://doi.org/10.1198/tast.2010.10052>
- Hui, F. K., Warton, D. I., & Foster, S. D. (2015). Tuning parameter selection for the adaptive lasso using ERIC. *Journal of the American Statistical Association*, 110, 262–269. <https://doi.org/10.1080/01621459.2014.951444>
- Hutchinson, R. A., Valente, J. J., Emerson, S. C., Betts, M. G., & Dietterich, T. G. (2015). Penalized likelihood methods improve parameter estimates in occupancy models. *Methods in Ecology and Evolution*, 6, 949–959. <https://doi.org/10.1111/2041-210X.12368>
- Johnston, A., Hochachka, W., Strimas-Mackey, M., Gutierrez, V. R., Robinson, O., Miller, E., ... Fink, D. (2019). Best practices for making reliable inferences from citizen science data: Case study using ebird to estimate species distributions. *bioRxiv*, 574392.
- Koshkina, V., Wang, Y., Gordon, A., Dorazio, R. M., White, M., & Stone, L. (2017). Integrated species distribution models: Combining presence-background data and site-occupancy data with imperfect detection. *Methods in Ecology and Evolution*, 8, 420–430.
- Kristensen, K., Nielsen, A., Berg, C. W., Skaug, H., & Bell, B. M. (2016). TMB: Automatic differentiation and Laplace approximation. *Journal of Statistical Software*, 70, 1–21. <https://doi.org/10.18637/jss.v070.i05>
- Lindgren, F., Rue, H., & Lindström, J. (2011). An explicit link between Gaussian fields and Gaussian Markov random fields: The stochastic partial differential equation approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 73, 423–498. <https://doi.org/10.1111/j.1467-9868.2011.00777.x>
- MacKenzie, D. I., Nichols, J. D., Hines, J. E., Knutson, M. G., & Franklin, A. B. (2003). Estimating site occupancy, colonization, and local extinction when a species is detected imperfectly. *Ecology*, 84, 2200–2207. <https://doi.org/10.1890/02-3090>
- Miller, D. A., Pacifici, K., Sanderlin, J. S., & Reich, B. J. (2019). The recent past and promising future for data integration methods to estimate species' distributions. *Methods in Ecology and Evolution*, 10(1), 22–37. <https://doi.org/10.1111/2041-210X.13110>
- Osborne, M. R., Presnell, B., & Turlach, B. A. (2000). On the lasso and its dual. *Journal of Computational and Graphical Statistics*, 9, 319–337.
- Pacifici, K., Reich, B. J., Miller, D. A., Gardner, B., Stauffer, G., Singh, S., ... Collazo, J. A. (2017). Integrating multiple data sources in species distribution modeling: A framework for data fusion. *Ecology*, 98, 840–850. <https://doi.org/10.1002/ecy.1710>
- Papangelou, F. (1974). The conditional intensity of general point processes and an application to line processes. *Probability Theory and Related Fields*, 28, 207–226.
- Phillips, S. J., Anderson, R. P., & Schapire, R. E. (2006). Maximum entropy modeling of species geographic distributions. *Ecological Modelling*, 190, 231–259. <https://doi.org/10.1016/j.ecolmodel.2005.03.026>
- Renner, I. W., Elith, J., Baddeley, A., Fithian, W., Hastie, T., Phillips, S. J., ... Warton, D. I. (2015). Point process models for presence-only analysis. *Methods in Ecology and Evolution*, 6, 366–379. <https://doi.org/10.1111/2041-210X.12352>
- Renner, I. W., & Warton, D. I. (2013). Equivalence of MAXENT and Poisson point process models for species distribution modeling in Ecology. *Biometrics*, 69, 274–281. <https://doi.org/10.1111/j.1541-0420.2012.01824.x>
- Simpson, D., Rue, H., Riebler, A., Martins, T. G., & Sørbye, S. H. (2017). Penalising model component complexity: A principled, practical approach to constructing priors. *Statistical Science*, 32, 1–28. <https://doi.org/10.1214/16-STS576>

- Sørbye, S. H., Illian, J. B., Simpson, D. P., Burslem, D., & Rue, H. (2019). Careful prior specification avoids incautious inference for log-gaussian cox point processes. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 68, 543–564. <https://doi.org/10.1111/rssc.12321>
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B (Methodological)*, 267–288. <https://doi.org/10.1111/j.2517-6161.1996.tb02080.x>
- Vandel, J. M., & Stahl, P. (2005). Distribution trend of the Eurasian lynx *Lynx lynx* populations in France. *Mammalia*, 69, 145–158. <https://doi.org/10.1515/mamm.2005.013>
- Warton, D. I., Renner, I. W., & Ramp, D. (2013). Model-based control of observer bias for the analysis of presence-only data in Ecology. *PLoS ONE*, 8, e79168. <https://doi.org/10.1371/journal.pone.0079168>
- Warton, D. I., & 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. <https://doi.org/10.1214/10-AOAS331>
- Widom, B., & Rowlinson, J. S. (1970). New model for the study of liquid-vapor phase transitions. *The Journal of Chemical Physics*, 52, 1670–1684. <https://doi.org/10.1063/1.1673203>
- Yuan, Y., Bachl, F. E., Lindgren, F., Borchers, D. L., Illian, J. B., Buckland, S. T., ... Gerrodette, T. (2017). Point process models for spatio-temporal distance sampling data from a large-scale survey of blue whales. *The Annals of Applied Statistics*, 11, 2270–2297. <https://doi.org/10.1214/17-AOAS1078>
- Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical Association*, 101, 1418–1429. <https://doi.org/10.1198/016214506000000735>

SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of the article.

How to cite this article: Renner IW, Louvrier J, Gimenez O. Combining multiple data sources in species distribution models while accounting for spatial dependence and overfitting with combined penalized likelihood maximization. *Methods Ecol Evol*. 2019;10:2118–2128. <https://doi.org/10.1111/2041-210X.13297>