

1 **Combining multiple data sources in species**  
2 **distribution models while accounting for**  
3 **spatial dependence and overfitting with**  
4 **combined penalised likelihood maximisation**

5 Ian W. Renner<sup>1\*</sup>, Julie Louvrier<sup>2</sup>, and Olivier Gimenez<sup>3</sup>

<sup>1</sup> School of Mathematical and Physical Sciences, The University of Newcastle, Australia

<sup>2</sup> Department of Ecological Dynamics, Leibniz Institute for Zoo and Wildlife Research,  
Department of Evolutionary Ecology, Alfred-Kowalke-Str. 17, D-10315 Berlin, Germany

<sup>3</sup> CEFE, CNRS, Univ Montpellier, Univ Paul Valéry Montpellier 3, EPHE, IRD, Montpellier, France

6 Word count: 6,785

\*ian.renner@newcastle.edu.au

## 7      **Summary**

- 8            1. The increase in availability of species data sets means that approaches  
9            to species distribution modelling that incorporate multiple data sets are  
10            in greater demand. Recent methodological developments in this area  
11            have led to combined likelihood approaches, in which a log-likelihood  
12            comprised of the sum of the log-likelihood components of each data source  
13            is maximised. Often, these approaches make use of at least one presence-  
14            only data set and use the log-likelihood of an inhomogeneous Poisson  
15            point process model in the combined likelihood construction. While these  
16            advancements have been shown to improve predictive performance, they  
17            do not currently address challenges in presence-only modelling such as  
18            checking and correcting for violations of the independence assumption  
19            of a Poisson point process model or more general challenges in species  
20            distribution modelling such as overfitting.
- 21            2. In this paper, we present an extension of the combined likelihood frame-  
22            work which accommodates alternative presence-only likelihoods in the  
23            presence of spatial dependence as well as lasso-type penalties to account  
24            for potential overfitting. We compare the proposed combined penalised  
25            likelihood approach to the standard combined likelihood approach via  
26            simulation and apply the method to modelling the distribution of the  
27            Eurasian lynx in the Jura Mountains in eastern France.
- 28            3. The simulations show that the proposed combined penalised likelihood  
29            approach has better predictive performance than the standard approach  
30            when spatial dependence is present in the data. The lynx analysis shows  
31            that the predicted maps vary significantly between the model fitted with  
32            the proposed combined penalised approach accounting for spatial depen-  
33            dence and the model fitted with the standard combined likelihood.
- 34            4. This work highlights the benefits of careful consideration of the presence-  
35            only components of the combined likelihood formulation, and allows  
36            greater flexibility and ability to accommodate real datasets.

37      **Keywords:** area-interaction models; diagnostic tools; lasso; occupancy models;  
38      point process models; presence-only data

## 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 *et al.*, 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 are available. Statistically, these methods are often fitted by maximising a corresponding likelihood expression, and the parameter estimates which maximise the likelihood may be used to produce maps of relative habitat suitability, reported as a habitat suitability index (Hirzel *et al.*, 2002), probability of species presence (Phillips *et al.*, 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 *et al.*, 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 *et al.* (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 maximising a combined log-likelihood expression which is the sum of the log-likelihoods of the presence-only and presence-absence components:

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

Here,  $\mathbf{s}_{\text{PO}}$  contains the locations of a presence-only data source, while  $\mathbf{y}_{\text{PA}}$  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 data sets are denoted by  $\boldsymbol{\alpha}_{\text{PO}}$  and  $\boldsymbol{\alpha}_{\text{PA}}$ , respectively, and collectively contained in the vector  $\boldsymbol{\alpha}$ . Hereafter, we refer to these parameters as sampling bias parameters, as

70 they may bias the intensity estimates as a result of the process of sampling the data. The  
71 key advancement of the combined likelihood approach is that the environmental response,  
72 parameterised by  $\beta$ , is informed by both the presence-only and presence-absence data.

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

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

80 While these papers clearly advance the practice of fitting SDMs in important ways, they  
81 do not address some common challenges that arise in real datasets. For example, they  
82 all consider an inhomogeneous Poisson point process model (IPPPM) for the presence-  
83 only data in the combination. In many real data sets, however, the implicit assumption  
84 that the point locations are independently distributed conditional on the environment is  
85 not met. Residual clustering or repulsion of the point locations not accounted for with  
86 an IPPPM due to the observation process, unconsidered environmental covariates, or  
87 biological factors would hence render the IPPPM inappropriate. One option to account  
88 for spatial dependence is to consider a log-Gaussian Cox Process, as Gelfand & Shirota  
89 (2018) do for a combination of presence-only and presence-absence data. Furthermore,  
90 none of the current literature in combined likelihood approaches includes ways to account  
91 for possible overfitting that results from including too many covariates in the model.

92 However, advances in SDM literature provide solutions to these common problems. Di-  
93 agnostic tools such as the inhomogeneous  $K$  function (Baddeley & Turner, 2000) and  
94 its simulation envelope (Diggle, 2003) can be used to determine departures from the  
95 independence assumption, and a wide number of alternative PPMs which account for  
96 spatial dependence may be included in the likelihood combination instead. Furthermore,  
97 penalised regression techniques such as the lasso penalty (Tibshirani, 1996) and its exten-  
98 sion the adaptive lasso (Zou, 2006) may be used as a way to perform variable selection.  
99 Lasso regularisation has been shown to boost predictive performance of SDMs and has  
100 been applied to IPPPMs (Renner & Warton, 2013) and occupancy models (Hutchinson  
101 *et al.*, 2015).

102 In this paper, we present a penalised combined likelihood model in a way that it is more

103 suitable for real data sets. In particular, we accommodate alternative forms of presence-  
104 only models to account for spatial dependence and affix a penalty on model complexity  
105 to address overfitting. In Section 2, we present the penalised combined likelihood formu-  
106 lation. In Section 3, we illustrate via simulations the improvements that this formulation  
107 provides and apply the proposed formulation to analyse the distribution of the Eurasian  
108 lynx (*lynx lynx*) in the Jura Mountains in eastern France. Finally, we present a discussion  
109 and further avenues for research in this area in Section 4.

## 110 2 Materials and Methods

### 111 2.1 Combined Penalised Likelihood Formulation

112 We define the weighted, combined penalised log-likelihood as follows

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

113 Here,  $\boldsymbol{\alpha} = (\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_D)^\top$  is a  $q$ -dimensional vector that collects coefficients for the vari-  
114 ables  $\mathbf{Z}$  used to model sampling bias for each of the  $D$  components individually. The  
115 environmental response is measured by a set of variables  $\mathbf{X}$  and is parametrised by  
116  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$ , which is collectively informed by all  $D$  components. The species  
117 data for all  $D$  components is collected in a set  $\mathbf{y}$ , with each individual data source  $\mathbf{y}_i$  de-  
118 termining the form of the component likelihood  $\ell_i(\boldsymbol{\alpha}_i, \boldsymbol{\beta}; \mathbf{y}_i)$ . Finally,  $p(\boldsymbol{\alpha}, \boldsymbol{\beta})$  is a penalty  
119 term described in further detail below.

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

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

126 In the simple occupancy model we consider, each site  $i$  is visited  $J_i$  times. We collect the  
127 history of detections and non-detections for all  $N$  sites in a matrix  $\mathbf{y}_{\text{occ}}$ . We assume that  
128 the probability that site  $i$  is occupied is given by  $\psi_i$  and that the occupancy of the sites  
129 remains constant throughout the history of visits. We further assume the probability of

130 detecting the species if present is  $p_i$ . Under these assumptions, we can then model the  
 131 probability of observing  $y_i$  detections at site  $i$  as

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

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

133 We can relate the occupancy  $\psi_i$  of site  $i$  to an inhomogeneous Poisson intensity  $\mu_i$  of the  
 134 species distribution over site  $i$  as in Koshkina *et al.* (2017):

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

135 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  
 136 is reasonable if  $\mu_i$  reasonably approximates the average intensity within site  $i$ .

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

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

142 The term  $p(\boldsymbol{\alpha}, \boldsymbol{\beta})$  in eqn 1 is a penalty on model complexity applied to both the envi-  
 143 ronmental parameters  $\boldsymbol{\beta}$  and the sampling bias parameters  $\boldsymbol{\alpha}$  to shrink these parameters  
 144 toward zero in order to boost predictive performance. Here, we consider both the tradi-  
 145 tional lasso penalty (Tibshirani, 1996) and the adaptive lasso penalty (Zou, 2006). For  
 146 the traditional lasso penalty,

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

147 where  $\lambda$  is the tuning parameter. For the adaptive lasso penalty,

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

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

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

149 for  $\gamma > 0$ . Here,  $\hat{\beta}_i^{(\text{unp})}$  is the unpenalised coefficient estimate corresponding to the  $i^{\text{th}}$  en-  
 150 vironmental variable  $\mathbf{x}_i$  and  $\hat{\alpha}_i^{(\text{unp})}$  is the unpenalised coefficient estimate corresponding to  
 151 the  $i^{\text{th}}$  sampling bias variable  $\mathbf{z}_i$ . The shape of the weights is determined by the parameter  
 152  $\gamma$ . The data-driven choice of the adaptive weights  $\mathbf{w}$  ensures that more important covari-  
 153 ates (*i.e.* those with coefficient estimates further away from 0) will be penalised less. This  
 154 construction also enables the adaptive lasso to achieve so-called oracle properties (Zou,  
 155 2006), which means that asymptotically, the correct subset of coefficients will be chosen  
 156 and the procedure has an optimal estimation rate.

157 We can use eqn 1 to represent the simpler framework introduced by Dorazio (2014) and  
 158 Fithian *et al.* (2015) by setting  $p(\boldsymbol{\alpha}, \boldsymbol{\beta}) = 0$ . We further extend this framework by  
 159 considering alternative choices for those component likelihoods  $\ell_i(\boldsymbol{\alpha}_i, \boldsymbol{\beta}; \mathbf{y}_i)$  informed by  
 160 presence-only data. Rather than consider only inhomogeneous Poisson point process  
 161 models, we consider area-interaction models (Widom & Rowlinson, 1970; Baddeley & van  
 162 Lieshout, 1995) when diagnostic analysis of these data sources identifies spatial depen-  
 163 dence among the presence-only locations. Area-interaction models account for spatial  
 164 dependence through a vector of computed point interactions  $\mathbf{t}_{\mathbf{s}}$ , which measure the pro-  
 165 portion of overlap among circles of a nominal radius around the observed points  $\mathbf{s}$ . They  
 166 can account for both clustering and repulsion of points – the model parameter  $\eta$  charac-  
 167 terises the nature of the spatial dependence, with values of  $\eta$  less than 1 signalling point  
 168 repulsion and values of  $\eta$  greater than 1 signalling point clustering.

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

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

171 This log-pseudolikelihood expression appears the same as eqn 2, with the exception that  
 172 the intensity  $\mu(s)$  is replaced by conditional intensity  $\mu(s; \mathbf{s}_{\text{PO}})$  (Papangelou, 1974), re-  
 173 flecting the fact that for the area-interaction model, intensity at a location  $s$  is conditional  
 174 on the other points in the pattern  $\mathbf{s}_{\text{PO}}$ .

## 175 2.2 Implementation in R

176 To fit models with the combined penalised log-likelihood in eqn 1, we have developed a set  
177 of functions in R inspired by the `optim` function and `ppmlasso` package (Renner & Warton,  
178 2013). The main function `comb_lasso` takes as an input a list of species data, associated  
179 environmental data, and formulae for the environmental trend and sampling bias trends  
180 for each component, along with details such as type of presence-only likelihoods to use,  
181 the type of penalty, the number of models to fit, and the tuning parameter criterion.  
182 The function applies the coordinate descent algorithm of Osborne *et al.* (2000). This  
183 requires the derivatives of the component likelihoods (also known as “score equations”)  
184 to be computed. Analytical score equations are supplied directly to the `optim` function,  
185 which serves as the machinery of the optimisation. A tutorial illustrating use of this code  
186 for the simulations as performed in Section 3.1 as well as some functions written to plot  
187 intensity maps and features of the lasso penalisation is provided in the supplementary  
188 material.

## 189 3 Results

### 190 3.1 Simulations

191 To investigate the benefits of the proposed penalised combined likelihood formulation, we  
192 used the `rpoispp` function in `spatstat` (Baddeley & Turner, 2005) to generate a large  
193 inhomogeneous Poisson pattern  $\mathbf{s}_{\text{true}}$  of roughly 10,000 points on a  $30 \times 30$ -unit square  
194 window from an intensity pattern defined by linear and quadratic terms of two gener-  
195 ated variables (hence four meaningful covariates  $\mathbf{x}_1, \dots, \mathbf{x}_4$  parameterised by coefficients  
196  $\beta_1, \dots, \beta_4$ ).

197 From this pattern, we generated two presence-only subsamples  $\mathbf{s}_1$  and  $\mathbf{s}_2$  biased by a  
198 different observation process. The first presence-only subsample  $\mathbf{s}_1$  was biased by  $\mathbf{z}_1$ , the  
199 distance to a simulated road network, and the other  $\mathbf{s}_2$  by  $\mathbf{z}_2$ , the distance to a simulated  
200 categorical covariate. We varied the size of the subsamples such that each pattern had  
201 25, 100, or 400 points. We also varied the strength of the clustering of the presence-only  
202 subsamples by setting the coefficient of the interaction term  $\nu_i = \ln \eta_i$  for  $i = 1, 2$ . Here,  
203 the patterns either exhibit no clustering ( $\nu_i = 0$ ), moderate clustering ( $\nu_i = 0.5$ ) or strong  
204 clustering ( $\nu_i = 1$ ). In each case, the radius of interactions is set to 1 spatial unit. To  
205 sample the points in  $\mathbf{s}_1$ , we proceed as follows:



- 206 1. Initialise the set of sampled points  $\mathbf{s}_1 = \emptyset$  and the point interactions  $\mathbf{t}_{\mathbf{s}_1}$  to be a  
207 vector of 0s
- 208 2. Compute the biased conditional intensity  $\mu_1(s; \mathbf{s}_1)$  at every point in  $\mathbf{s}_{\text{true}}$  using  
209  $\mathbf{x}_1, \dots, \mathbf{x}_4$ , the sampling bias covariate  $\mathbf{z}_1$ , and the current vector of point inter-  
210 actions  $\mathbf{t}_{\mathbf{s}_1}$ , where the biased conditional intensity is defined as follows:  
$$\ln \mu_1(s; \mathbf{s}_1) = \beta_1 \mathbf{x}_1(s) + \beta_2 \mathbf{x}_2(s) + \beta_3 \mathbf{x}_3(s) + \beta_4 \mathbf{x}_4(s) + \alpha_1 \mathbf{z}_1(s) + \nu_1 \mathbf{t}_{\mathbf{s}_1}(s)$$
- 211 3. Set  $\mu_1(s; \mathbf{s}_1) = 0$  for all  $s \in \mathbf{s}_1$ . That is, we set the conditional intensity for any  
212 point already selected in  $\mathbf{s}_1$  to 0 to ensure these points are not resampled
- 213 4. Randomly select a point from  $\mathbf{s}_{\text{true}}$  with sampling probabilities proportional to the  
214 conditional intensities and add the selected point to  $\mathbf{s}_1$
- 215 5. Update the vector of point interactions  $\mathbf{t}_{\mathbf{s}_1}$  for all points in  $\mathbf{s}_{\text{true}}$  using the internal  
216 `evalInteraction` function in `spatstat`, which computes point interactions based  
217 on a supplied point pattern for a given set of locations and interaction radius
- 218 6. Repeat steps 2-5 until we have sampled the desired number of points

219 We sample  $\mathbf{s}_2$  in a similar manner, using  $\mathbf{z}_2$  instead of  $\mathbf{z}_1$  to create the sampling bias and  
220 computing point interactions  $\mathbf{t}_{\mathbf{s}_2}$ .

221 Because the true pattern  $\mathbf{s}_{\text{true}}$  is Poisson, this simulation setup emulates a scenario in which  
222 the clustering of the observed point patterns is an artefact of the observation process –  
223 this can happen if, for example, records are publicly available and enthusiasts for the  
224 species report further observations near the publicly available locations (Johnston *et al.*,  
225 2019).

226 We also generated a history  $\mathbf{y}_{\text{occ}}$  of detections and non-detections from 5 visits to each of  
227 100 sites centred along a regular grid in the  $30 \times 30$ -unit observation window to emulate  
228 a data set for which we could consider occupancy modelling. The species was considered  
229 present at a site if the closest point in the pattern  $\mathbf{s}_{\text{true}}$  was within a distance of 0.18  
230 units of the centre of the site, such that the area of each site is roughly 0.1 square units.  
231 The history of detections and non-detections at each site where the species was considered  
232 present was randomly generated according to detection probabilities defined by the inverse  
233 of the cloglog function evaluated at a generated detection covariate  $\mathbf{z}_3$ .

234 Finally, we generated four dummy covariates  $\mathbf{d}_1, \dots, \mathbf{d}_4$  to include in fitted models that  
 235 were meaningless in describing the true species distribution. We did this to reflect the  
 236 fact that in real applications, we may not know which among a suite of candidate vari-  
 237 ables truly determine the species distribution. We ensured that the maximum absolute  
 238 correlation among all pairs of variables was smaller than 0.5.

239 After generating the species data, we fit a number of models, using as input environmental  
 240 covariates the four meaningful covariates  $\mathbf{x}_1, \dots, \mathbf{x}_4$  (parameterised by  $\beta_1, \dots, \beta_4$ ) as well  
 241 as four dummy covariates  $\mathbf{d}_1, \dots, \mathbf{d}_4$  (parameterised by  $\beta_5, \dots, \beta_8$ ) and using as sampling  
 242 bias covariates  $\mathbf{z}_1, \mathbf{z}_2$ , and  $\mathbf{z}_3$  (parameterised by  $\alpha_1, \alpha_2$ , and  $\alpha_3$ ). For both Poisson and  
 243 area-interaction presence-only likelihoods, we fit a model without any penalty, with a  
 244 lasso penalty, and with an adaptive lasso penalty. For the models fitted with either a  
 245 lasso or an adaptive lasso penalty, we fit regularisation paths of 1000 models, increasing  
 246 the penalty from 0 to the smallest penalty  $\lambda_{\max}$  that would shrink all coefficients to 0,  
 247 thus covering the entire scope of possible model sizes. The model which minimised BIC  
 248 was chosen among the 1000 fitted models. We considered as species data a combination  
 249 of all three of  $\mathbf{s}_1, \mathbf{s}_2$ , and  $\mathbf{y}_{\text{occ}}$ . This led to a total of six models being fitted, summarised  
 250 in Table 1.

Model	Species Data	Presence-only likelihood	Penalty
1	$\mathbf{s}_1, \mathbf{s}_2$ , and $\mathbf{y}_{\text{occ}}$	IPPPM	None
2	$\mathbf{s}_1, \mathbf{s}_2$ , and $\mathbf{y}_{\text{occ}}$	IPPPM	Lasso
3	$\mathbf{s}_1, \mathbf{s}_2$ , and $\mathbf{y}_{\text{occ}}$	IPPPM	Adaptive Lasso
4	$\mathbf{s}_1, \mathbf{s}_2$ , and $\mathbf{y}_{\text{occ}}$	Area-interaction	None
5	$\mathbf{s}_1, \mathbf{s}_2$ , and $\mathbf{y}_{\text{occ}}$	Area-interaction	Lasso
6	$\mathbf{s}_1, \mathbf{s}_2$ , and $\mathbf{y}_{\text{occ}}$	Area-interaction	Adaptive Lasso

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

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

256 We performed 1,000 simulations of the data sets for each of the nine combinations of

257 presence-only data set size and clustering strength and the resultant model fits on 512GB  
258 nodes powered by 3.0 GHz Intel Xeon Gold (E5-6154) processor from the University  
259 of Newcastle’s High Performance Computing cluster. The 9,000 simulation tasks took  
260 approximately 7,000 hours.

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

271 In the Appendix, we show that the parameter coefficients  $\beta_1, \dots, \beta_4$  corresponding to the  
272 meaningful covariates  $\mathbf{x}_1, \dots, \mathbf{x}_4$  are increasingly biased away from 0 for the models fitted  
273 with Poisson presence-only likelihoods, both as sample size increases and as the strength  
274 of presence-only clustering increases. The inclusion of the area-interaction term takes  
275 an increasing amount of signal from the environmental covariates as the strength of the  
276 presence-only clustering increases. For low sample sizes, there is a suggestion that this  
277 signal dampening may be too strong, though such an overcorrection disappears as sample  
278 size increases.

279 Second, penalisation via the lasso or adaptive lasso improves model performance when  
280 there is no presence-only clustering, and this improvement is greatest for smaller sam-  
281 ple sizes. This is an expected conclusion given the danger of overfitting is greater with  
282 fewer observations. Models penalised with the adaptive lasso tend to outperform models  
283 penalised with the lasso when there is no presence-only clustering. However, lasso penal-  
284 isation does not notably improve performance when there is presence-only clustering. In  
285 fact, there is a suggestion that applying a lasso penalty may slightly hinder performance  
286 when applying an area-interaction presence-only likelihood for small sample sizes. Al-  
287 though the benefits of penalisation are negligible with large data sets, fitting models with  
288 a penalty does not hurt the performance.

289 In summary, it appears that the proposed combined penalised likelihood framework pro-  
290 vides the best performance. Furthermore, incorporating area-interaction presence-only

291 likelihoods improves performance when clustering is present, and can likewise reliably  
 292 estimate that there are negligible point interactions if clustering is not present, in effect  
 293 relaxing to the simpler model with Poisson presence-only likelihoods when this additional  
 294 complexity is not needed. A more detailed discussion of the simulation results, including  
 295 boxplots of the fitted coefficients, appears in the Appendix.

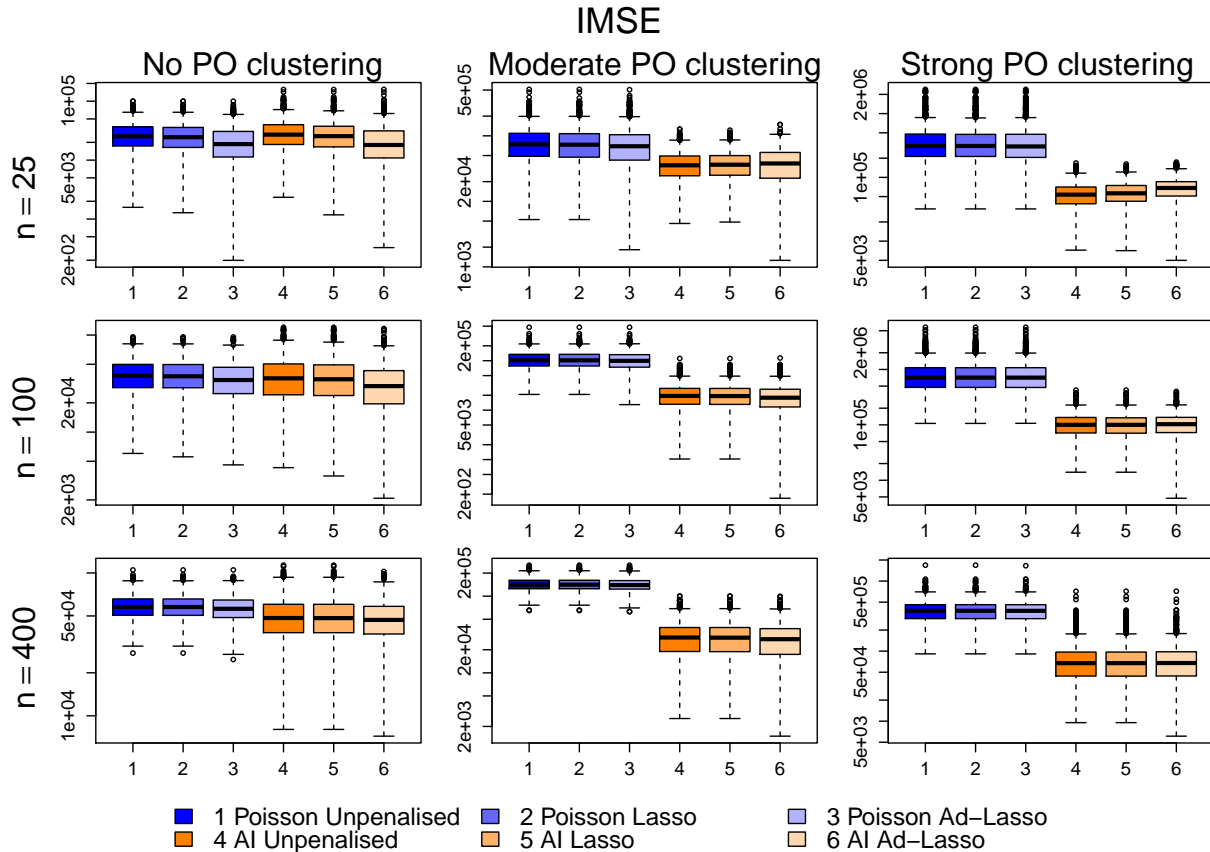


Figure 1: Boxplots of integrated mean squared error for the six models described in Table 1 for different combinations of presence-only sample size and clustering strength.

### 296 **3.2 Analysis of Eurasian lynx distribution in the Jura Moun-** 297 **tains**

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

300 Lynx went extinct in France at the end of the 19th century due to habitat degradation,  
 301 human persecution and decrease in prey availability (Vandel & Stahl, 2005). The species  
 302 was reintroduced in Switzerland in the 1970s (Breitenmoser *et al.*, 1998), then re-colonised  
 303 France through the Jura mountains in the 1980s (Vandel & Stahl, 2005). The species is

304 listed as endangered under the 2017 IUCN Red list and is of conservation concern in  
305 France due to habitat fragmentation, poaching and collisions with vehicles. The Jura  
306 holds the bulk of the French lynx population.

307 We have three sources of lynx data in the Jura Mountains: a presence-only data set  
308 consisting of 440 opportunistic sightings in the wild from 2009-2011 (denoted  $\mathbf{s}_w$ ), another  
309 presence-only data set consisting of 240 reported interferences of lynx with domestic  
310 livestock in 2009-2011 (denoted  $\mathbf{s}_d$ ), and pictures of lynx taken from cameras set up in  
311 73 locations  $\mathbf{s}_c$  in the Jura Mountains in 2012. Lynx presence-only data were made  
312 of presence signs sampled all year long thanks to a network of professional and non-  
313 professional observers. Every observer is trained during a 3-day teaching course led by the  
314 French National Game and Wildlife Agency (ONCFS) to document signs of the species'  
315 presence (Duchamp *et al.*, 2012). Presence signs went through a standardised control  
316 process to prevent misidentification (Duchamp *et al.*, 2012). The camera data has daily  
317 reportings of the lynx across a total of 77 days. Due to this, we can consider the picture  
318 history of lynx at the camera locations in an occupancy modelling framework (Blanc *et al.*,  
319 2014). In particular, we split the 77-day period into seven 11-day periods, such that the  
320 site history  $\mathbf{y}_c$  comprises seven detections and non-detections at each site in  $\mathbf{s}_c$  over each  
321 11-day period.

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

328 To model the lynx distribution, we consider altitude, percentage of forest cover, distance  
329 to the nearest water source, and human population density as environmental variables.  
330 We model sampling bias in the wild records  $\mathbf{s}_w$  with distance to the nearest main road  
331 and distance to the nearest train line, and sampling bias in the domestic records  $\mathbf{s}_d$   
332 with distance to the nearest farm and percentage of agricultural land. Finally, we model  
333 detection probability for the camera data with distance to the nearest urban area. We  
334 established this set of potential candidate environmental and detection variables based on  
335 previously studied species habitat preferences and detectability (Bouyer *et al.*, 2015). The  
336 Corine Land Cover land use repository from 2012 (Büttner *et al.*, 2014) supplies a map of  
337 land coverage including urban areas, water areas, forest areas, farm areas, and agricultural

## Lynx Locations

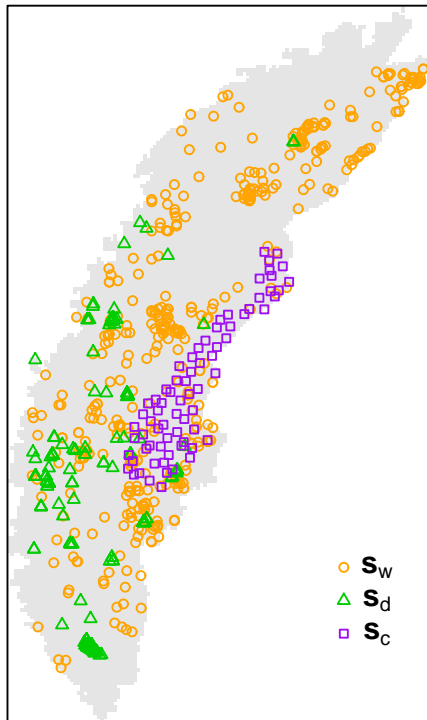


Figure 2: Locations of the lynx data in the Jura Mountains, including 440 observations in the wild  $s_w$ , 220 reports of domestic interference  $s_d$ , and 73 camera traps  $s_c$ .

338 areas that was used to generate the percentage of forest areas and agricultural areas over  
339  $1 \text{ km} \times 1 \text{ km}$  cells as well as distances to the nearest urban area, water source, and farm.  
340 Altitude was averaged over  $1 \text{ km} \times 1 \text{ km}$  cells from data available in the `raster` package  
341 in R, while human population density was averaged over  $1 \text{ km} \times 1 \text{ km}$  cells taken from  
342 version 4 of the Gridded Population of the World data repository (Center for International  
343 Earth Science Information Network (CIESIN) – Columbia University, 2016). Distances  
344 from the nearest main road and railway were computed from shapefiles from Version 151  
345 of the ROUTE 500 database, accessible at <http://professionnels.ign.fr/route500>.

346 We fitted initial separate IPPPMs to the wild records  $s_w$  and the domestic records  $s_d$  using  
347 linear, quadratic, and interaction terms for the four environmental covariates, and linear  
348 terms for the sampling bias covariates. From these models, we are able to assess whether  
349 the assumption of independence inherent to the IPPPMs is appropriate with simulation  
350 envelopes of the inhomogeneous  $K$ -function in `spatstat`, as shown in Figure 3. Both  
351 of the envelopes for the IPPPMs fitted to the wild model (left panel) and the domestic  
352 model (middle panel) demonstrate additional clustering as the observed inhomogeneous  
353  $K$ -function values plotted in red fall above the simulation envelopes for small radii. This

354 suggests that fitting an IPPPM is inappropriate for these data sets. The right panel shows  
355 a simulation envelope of the cross  $K$  function as produced by the `Kcross.inhom` function  
356 of `spatstat`, which counts the expected number of wild sightings within a given distance  
357 of a domestic sighting, conditional on the spatially varying intensities of both patterns.  
358 We estimate the wild and domestic intensities from area-interaction models, and as the  
359 observed values of the cross  $K$ -function fall within the envelope boundaries, this suggests  
360 that there is no clustering across the two data sets. This, in turn, suggests that the  
361 observed clustering within the wild and domestic data sets may be more likely attributable  
362 to the observation process than to some biological reality that induces clustering or a  
363 missed environmental covariate.

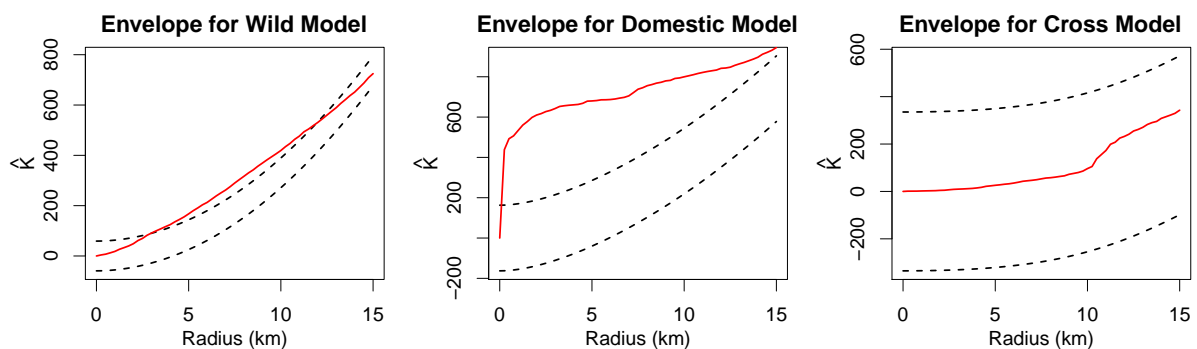


Figure 3: 95% simulation envelopes of the inhomogeneous  $K$ -function for the fitted IPPPM of the wild records (left), the fitted IPPPM of the domestic records (middle), and across the two fitted IPPPMs (right).

364 Consequently, we fit combined likelihood models using both the standard, unpenalised  
365 approach (analogous to Model 1 in Table 1) and the combined penalised likelihood for-  
366 mulation eqn 1 with a lasso penalty and area-interaction models for the presence-only  
367 data sources (analogous to Model 5 in Table 1). The radii chosen to capture the residual  
368 spatial patterning in the wild and domestic models are 2km and 5km, as chosen by the  
369 `profilepl` function in `spatstat`.

370 Figure 4 shows the bias-corrected fitted intensities from these two models. For the com-  
371 bined model which uses IPPPMs (left panel), the fitted intensity is corrected for the sam-  
372 pling bias terms modelled for the presence-only components using the method of Warton  
373 *et al.* (2013). For the combined penalised model which uses area-interaction models (right  
374 panel), the fitted intensity is corrected for these same sampling bias terms as well as the  
375 fitted point interactions – that is, we treat the interaction parameter  $\nu$  as belonging to the  
376 set of sampling bias parameters  $\alpha$ . The fitted models show strikingly different patterns,  
377 with the model which uses area-interaction components highlighting much more of the

378 Jura Mountains as preferred habitat of lynx than the model which uses IPPPMs. The  
379 models suggest similar numbers of points throughout the Jura, but the distribution of  
380 these points are more heavily concentrated in the IPPPM model. This is because the  
381 area-interaction terms in the AI model lessen the impact of some clusters of points on the  
382 scale of the displayed bias-corrected intensities.

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

387 The combined penalised model with the area-interaction components found the optimal  
388 lasso penalty was 0, resulting in a model which included all 18 covariates and both of  
389 the area-interaction terms. The fact that the optimal penalty is 0 suggests that the suite  
390 of covariates we chose to include, motivated by existing literature, seems to have been a  
391 good choice. In general, we recommend use of the lasso penalty as a safeguard against  
392 overfitting, particularly in contexts where the suite of candidate covariates for a species  
393 is less established as an insurance against overfitting.

## 394 4 Discussion

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

### 399 4.1 Possible extensions

400 Despite these improvements, further advances are possible. Other penalty structures  
401 could be incorporated into the same framework. While the lasso and adaptive lasso  
402 showcased here show clear benefits in simulations, other penalised likelihood variants  
403 such as SCAD (Fan & Li, 2001) could lead to superior performance in some situations,  
404 and alternative methods to BIC of choosing the size of the penalty such as the Extended  
405 Bayesian Information Criterion (“ERIC”, Hui *et al.*, 2015) could likewise be used.

406 While we make use of the area-interaction likelihood in this paper, there is a large family of  
407 Gibbs PPMs (Cressie, 1992) which accommodate different sorts of spatial dependence that



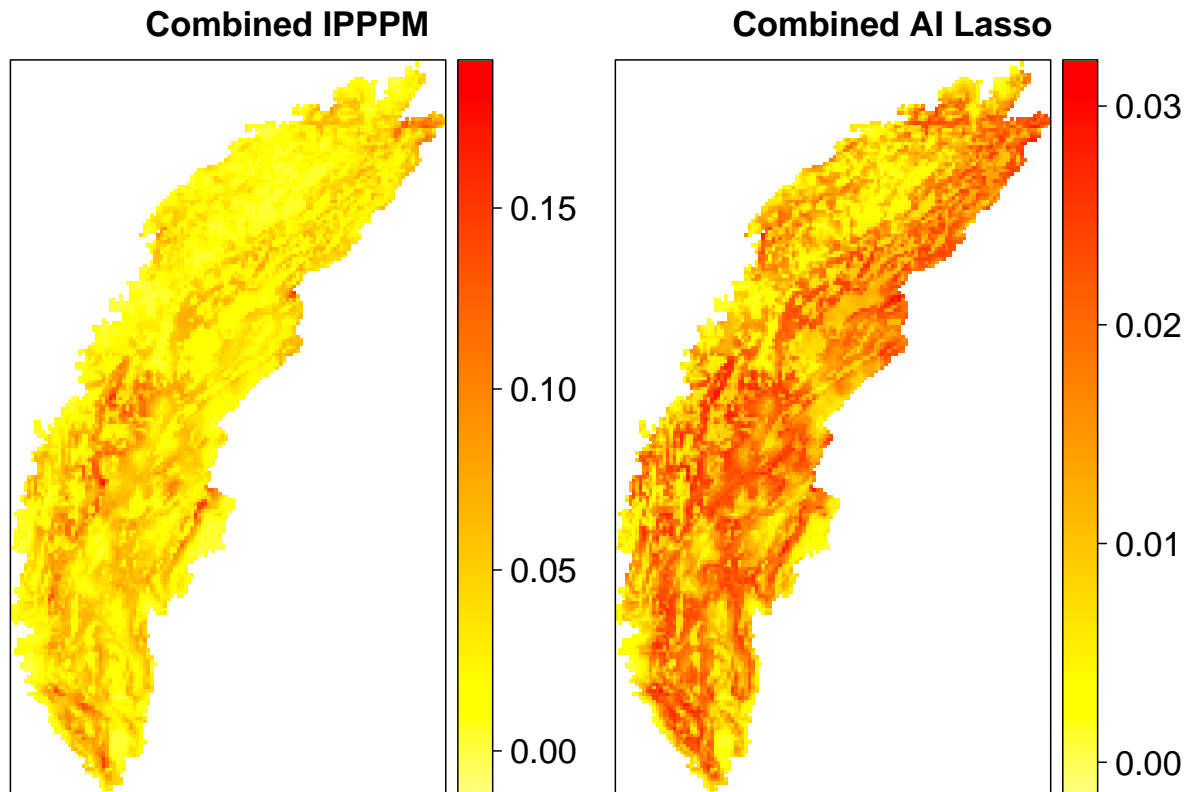


Figure 4: Fitted intensities using the combined likelihood formulation. Left: the model is fitted without any penalty and using inhomogeneous Poisson point process models for the presence-only data sources. Right: the model is fitted with a lasso penalty and using area-interaction models for the presence-only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

408 could be used. Our choice of the area-interaction model as the alternative is motivated by  
409 the fact that it accommodates interactions of all orders instead of just pairwise interactions  
410 and that it can be used to model both clustering and repulsion of points.

411 The inclusion of the area-interaction terms dampens the signal of the environmental co-  
412 variates. Although this makes sense when spatial dependence exists, we may dampen the  
413 signal too much. In the context of species distribution models, we might ask the question,  
414 “Does a given species record exist because its location is in particularly suitable habitat  
415 for the species, or because there are other records nearby?” If the answer to this question  
416 appears to be “both”, as is often the case for presence-only data, we are at risk of “spatial  
417 confounding”. In the single-source context, Hodges & Reich (2010) propose restricting the  
418 spatial effect to be orthogonal to the fixed covariate effects, while Simpson *et al.* (2017)  
419 and Sørbye *et al.* (2019) suggest careful selection of associated spatial priors to alleviate  
420 this risk. With our implementation, we could achieve something similar to the latter two  
421 papers by adjusting the magnitude of the lasso penalty for the area-interaction terms. In

422 the Appendix, we highlight the tradeoff between the estimates of the interaction parame-  
423 ters  $\hat{\nu}_i$  and both the estimates of the environmental parameters  $\hat{\beta}_i$  and the sampling bias  
424 parameters  $\hat{\alpha}_i$ . However, a full exploration of the effects of spatial confounding remains  
425 an open area of research and is beyond the scope of this paper.

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

434 Further improvements could be made by incorporating source weights in situations in  
435 which the data sources vary in quality. Indeed, presence-only data sources may be more  
436 prone to errors in coordinate locations as well as correct species identification, as they often  
437 include records by amateur enthusiasts. The combined penalised likelihood framework  
438 could easily be extended to include weights for the various data sources by adding a  
439 vector of source weights  $\mathbf{w} = (w_1, \dots, w_D)^\top$  to the formulation in eqn 1:

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

440 One possible strategy to incorporate such weights in eqn 3 could be to compare perfor-  
441 mance of single source models on independent data and upweight the contribution of data  
442 sources that are shown to have good performance.

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

## 446 4.2 Accounting for dependence within and among data sources

447 In the lynx data analysis in Section 3.2, we diagnosed spatial dependence within each  
448 of the presence-only data sources but found no spatial dependence across data sources.  
449 Tools such as the inhomogeneous  $K$ -envelope provide great insight into the underlying  
450 individual spatial processes that are observed. However, such diagnostic tools are not

451 currently available for the combined likelihood models, and research in this area would  
452 be valuable as these models grow in popularity.

453 Another approach to constructing SDMs from multiple data sources could be to introduce  
454 a common latent spatial term  $\xi(s)$ , such as a Gaussian random field, which would account  
455 for spatial dependence among points in all of the data sources as in Gelfand & Shirota  
456 (2018). The resulting likelihood expression would be:

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

457 where  $\xi(\mathbf{y}) \sim \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma})$ . Models of this type are typically fitted in a Bayesian frame-  
458 work. We could reduce the dimension of  $\xi$  through methods like fixed rank kriging or  
459 induce sparsity in  $\boldsymbol{\Sigma}$  through lasso-type penalties such that the likelihood in eqn 4 could  
460 be fitted with software such as Template Model Builder (TMB, Kristensen *et al.*, 2016).  
461 Another way to achieve sparsity is with the stochastic partial differential equation ap-  
462 proach (SPDE, Lindgren *et al.*, 2011), as implemented in the `inlabru` package (Bachl  
463 *et al.*, 2019).

### 464 4.3 Conclusion and Perspectives

465 The development of statistical methods is often motivated by new challenges raised by  
466 novel types of data sets. While the current literature on combined likelihood approaches  
467 represents a significant recent advancement, advances in other areas can be lost if not  
468 carried over with such methodological developments. This paper attempts to build a  
469 bridge between this exciting new arena for species distribution modelling and the rich  
470 suite of tools available for species distribution modelling, particularly that for presence-  
471 only data. Our hope is that other such bridges continue to be built in this spirit.

## 472 5 Acknowledgements

473 We thank the staff from the French National Game and Wildlife Agency, the Forest  
474 National Agency and the Departmental Federation of Hunters of Jura department, who  
475 collected the photographs during the camera-trapping session. We also thank all the  
476 volunteers that are members of the Réseau Loup-Lynx that collect every year precious  
477 presence signs of lynx all over its distribution area. We thank the University of Newcastle

478 through the Early to Mid-Career Researcher Visiting Fellowship and the Université Paul-  
479 Valéry Montpellier through the Professeurs en mobilité universitaire fund for funding  
480 visits for I.R. which helped facilitate the collaborations that resulted in this work. O.G.  
481 was funded by CNRS and the “Mission pour l’Interdisciplinarité” through the “Osez  
482 l’Interdisciplinarité” initiative. Finally, we thank the reviewers for their very helpful  
483 comments, which have helped us greatly improve the paper.

## 484 6 Authors’ Contributions

485 I.R. and O.G. conceived the concept of the paper. I.R. developed the code to fit the  
486 models. J.L. sourced the species coordinates and covariates for the lynx analysis. I.R.  
487 and O.G. wrote the manuscript. I.R. and J.L. developed the tutorial in the supplementary  
488 information. All authors were involved in editing drafts of the manuscript.

## 489 References

- 490 Bachl, F.E., Lindgren, F., Borchers, D.L. & Illian, J.B. (2019) inlabru: an r package  
491 for bayesian spatial modelling from ecological survey data. *Methods in Ecology and*  
492 *Evolution*, **10**, 760–766.
- 493 Baddeley, A. & Turner, R. (2000) Practical maximum pseudolikelihood for spatial point  
494 patterns: (with discussion). *Australian & New Zealand Journal of Statistics*, **42**, 283–  
495 322.
- 496 Baddeley, A. & Turner, R. (2005) Spatstat: an R package for analyzing spatial point  
497 patterns. *Journal of Statistical Software*, **12**, 1–42.
- 498 Baddeley, A.J. & van Lieshout, M.N.M. (1995) Area-interaction point processes. *Annals*  
499 *of the Institute of Statistical Mathematics*, **47**, 601–619.
- 500 Besag, J. (1977) Some methods of statistical analysis for spatial data. *Bulletin of the*  
501 *International Statistical Institute*, **47**, 77–92.
- 502 Blanc, L., Marboutin, E., Gatti, S., Zimmermann, F. & Gimenez, O. (2014) Improving  
503 abundance estimation by combining capture–recapture and occupancy data: example  
504 with a large carnivore. *Journal of Applied Ecology*, **51**, 1733–1739.

- 505 Bouyer, Y., San Martin, G., Poncin, P., Beudels-Jamar, R.C., Odden, J. & Linnell, J.D.  
506 (2015) Eurasian lynx habitat selection in human-modified landscape in norway: Effects  
507 of different human habitat modifications and behavioral states. *Biological Conservation*,  
508 **191**, 291–299.
- 509 Breitenmoser, U., Breitenmoser-Würsten, C. & Capt, S. (1998) Re-introduction and  
510 present status of the lynx (*lynx lynx*) in switzerland. *Hystrix, the Italian Journal*  
511 *of Mammalogy*, **10**.
- 512 Büttner, G., Soukup, T. & Kosztra, B. (2014) Clc2012 addendum to clc2006 technical  
513 guidelines. *Final Draft, Copenhagen (EEA)*.
- 514 Center for International Earth Science Information Network (CIESIN) – Columbia Uni-  
515 versity (2016) Gridded population of the world, version 4 (gpwv4): population density.
- 516 Cressie, N. (1992) *Statistics for spatial data*, volume 4. Wiley Online Library.
- 517 Cressie, N. & Wikle, C.K. (2015) *Statistics for spatio-temporal data*. John Wiley & Sons.
- 518 Diggle, P.J. (2003) *Statistical analysis of spatial point patterns*. Edward Arnold.
- 519 Dorazio, R.M. (2014) Accounting for imperfect detection and survey bias in statistical  
520 analysis of presence-only data. *Global Ecology and Biogeography*, **23**, 1472–1484.
- 521 Duchamp, C., Boyer, J., Briaudet, P.E., Leonard, Y., Moris, P., Bataille, A., Dahier, T.,  
522 Delacour, G., Millisher, G., Miquel, C. *et al.* (2012) A dual frame survey to assess time-  
523 and space-related changes of the colonizing wolf population in france. *Hystrix-Italian*  
524 *Journal of Mammalogy*, **23**, 14–28.
- 525 Fan, J. & Li, R. (2001) Variable selection via nonconcave penalized likelihood and its  
526 oracle properties. *Journal of the American statistical Association*, **96**, 1348–1360.
- 527 Fithian, W., Elith, J., Hastie, T. & Keith, D.A. (2015) Bias correction in species dis-  
528 tribution models: pooling survey and collection data for multiple species. *Methods in*  
529 *Ecology and Evolution*, **6**, 424–438.
- 530 Gelfand, A. & Shirota, S. (2018) Preferential sampling for presence/absence data  
531 and for fusion of presence/absence data with presence-only data. *arXiv preprint*  
532 *arXiv:180901322*.

- 533 Gould, M.J., Gould, W.R., Cain III, J.W. & Roemer, G.W. (2019) Validating the perfor-  
534 mance of occupancy models for estimating habitat use and predicting the distribution  
535 of highly-mobile species: A case study using the american black bear. *Biological Con-*  
536 *servation*, **234**, 28–36.
- 537 Guisan, A., Thuiller, W. & Zimmermann, N.E. (2017) *Habitat suitability and distribution*  
538 *models: with applications in R*. Cambridge University Press.
- 539 Hirzel, A.H., Hausser, J., Chessel, D. & Perrin, N. (2002) Ecological-niche factor analysis:  
540 how to compute habitat-suitability maps without absence data? *Ecology*, **83**, 2027–  
541 2036.
- 542 Hodges, J.S. & Reich, B.J. (2010) Adding spatially-correlated errors can mess up the fixed  
543 effect you love. *The American Statistician*, **64**, 325–334.
- 544 Hui, F.K., Warton, D.I. & Foster, S.D. (2015) Tuning parameter selection for the adaptive  
545 lasso using ERIC. *Journal of the American Statistical Association*, **110**, 262–269.
- 546 Hutchinson, R.A., Valente, J.J., Emerson, S.C., Betts, M.G. & Dietterich, T.G. (2015) Pe-  
547 nalized likelihood methods improve parameter estimates in occupancy models. *Methods*  
548 *in Ecology and Evolution*, **6**, 949–959.
- 549 Johnston, A., Hochachka, W., Strimas-Mackey, M., Gutierrez, V.R., Robinson, O., Miller,  
550 E., Auer, T., Kelling, S. & Fink, D. (2019) Best practices for making reliable inferences  
551 from citizen science data: case study using ebird to estimate species distributions.  
552 *bioRxiv*, p. 574392.
- 553 Koshkina, V., Wang, Y., Gordon, A., Dorazio, R.M., White, M. & Stone, L. (2017)  
554 Integrated species distribution models: combining presence-background data and site-  
555 occupancy data with imperfect detection. *Methods in Ecology and Evolution*, **8**, 420–430.
- 556 Kristensen, K., Nielsen, A., Berg, C.W., Skaug, H. & Bell, B.M. (2016) TMB: Automatic  
557 differentiation and Laplace approximation. *Journal of Statistical Software*, **70**, 1–21.
- 558 Lindgren, F., Rue, H. & Lindström, J. (2011) An explicit link between gaussian fields and  
559 gaussian markov random fields: the stochastic partial differential equation approach.  
560 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **73**, 423–  
561 498.

- 562 MacKenzie, D.I., Nichols, J.D., Hines, J.E., Knutson, M.G. & Franklin, A.B. (2003)  
563 Estimating site occupancy, colonization, and local extinction when a species is detected  
564 imperfectly. *Ecology*, **84**, 2200–2207.
- 565 Miller, D.A., Pacifici, K., Sanderlin, J.S. & Reich, B.J. (2019) The recent past and promis-  
566 ing future for data integration methods to estimate species' distributions. *Methods in*  
567 *Ecology and Evolution*, **10**, 22–37.
- 568 Osborne, M.R., Presnell, B. & Turlach, B.A. (2000) On the lasso and its dual. *Journal*  
569 *of Computational and Graphical Statistics*, **9**, 319–337.
- 570 Pacifici, K., Reich, B.J., Miller, D.A., Gardner, B., Stauffer, G., Singh, S., McKerrow,  
571 A. & Collazo, J.A. (2017) Integrating multiple data sources in species distribution  
572 modeling: A framework for data fusion. *Ecology*, **98**, 840–850.
- 573 Papangelou, F. (1974) The conditional intensity of general point processes and an appli-  
574 cation to line processes. *Probability Theory and Related Fields*, **28**, 207–226.
- 575 Phillips, S.J., Anderson, R.P. & Schapire, R.E. (2006) Maximum entropy modeling of  
576 species geographic distributions. *Ecological Modelling*, **190**, 231–259.
- 577 Renner, I.W. & Warton, D.I. (2013) Equivalence of MAXENT and Poisson point process  
578 models for species distribution modeling in ecology. *Biometrics*, **69**, 274–281.
- 579 Renner, I.W., Elith, J., Baddeley, A., Fithian, W., Hastie, T., Phillips, S.J., Popovic,  
580 G. & Warton, D.I. (2015) Point process models for presence-only analysis. *Methods in*  
581 *Ecology and Evolution*, **6**, 366–379.
- 582 Simpson, D., Rue, H., Riebler, A., Martins, T.G., Sørbye, S.H. *et al.* (2017) Penalising  
583 model component complexity: A principled, practical approach to constructing priors.  
584 *Statistical Science*, **32**, 1–28.
- 585 Sørbye, S.H., Illian, J.B., Simpson, D.P., Burslem, D. & Rue, H. (2019) Careful prior  
586 specification avoids incautious inference for log-gaussian cox point processes. *Journal*  
587 *of the Royal Statistical Society: Series C (Applied Statistics)*, **68**, 543–564.
- 588 Tibshirani, R. (1996) Regression shrinkage and selection via the lasso. *Journal of the*  
589 *Royal Statistical Society Series B (Methodological)*, pp. 267–288.
- 590 Vandel, J.M. & Stahl, P. (2005) Distribution trend of the eurasian lynx lynx lynx popu-  
591 lations in france. *Mammalia mamm*, **69**, 145–158.

- 592 Warton, D.I., Renner, I.W. & Ramp, D. (2013) Model-based control of observer bias for  
593 the analysis of presence-only data in ecology. *PloS one*, **8**, e79168.
- 594 Warton, D.I. & Shepherd, L.C. (2010) Poisson point process models solve the “pseudo-  
595 absence problem” for presence-only data in ecology. *Annals of Applied Statistics*, **4**,  
596 1383–1402.
- 597 Widom, B. & Rowlinson, J.S. (1970) New model for the study of liquid–vapor phase  
598 transitions. *The Journal of Chemical Physics*, **52**, 1670–1684.
- 599 Yuan, Y., Bachl, F.E., Lindgren, F., Borchers, D.L., Illian, J.B., Buckland, S.T., Rue, H.,  
600 Gerrodette, T. *et al.* (2017) Point process models for spatio-temporal distance sampling  
601 data from a large-scale survey of blue whales. *The Annals of Applied Statistics*, **11**,  
602 2270–2297.
- 603 Zou, H. (2006) The adaptive lasso and its oracle properties. *Journal of the American*  
604 *Statistical Association*, **101**, 1418–1429.