

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^{1*}, Julie Louvrier², and Olivier Gimenez³

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

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

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

6 Word count: 6,060

*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 approach combined
25 penalised likelihood approach to the standard combined likelihood ap-
26 proach via simulation and apply the method to modelling the distribution
27 of the Eurasian lynx in the Jura Mountains in eastern France.
- 28 3. The simulations show that the proposed combined penalised likelihood
29 approach outperforms the standard approach when spatial dependence is
30 present in the data. The lynx analysis shows that the predicted maps vary
31 significantly with the different implementations of the proposed approach.
- 32 4. This work highlights the benefits of careful consideration of the presence-
33 only components of the combined likelihood formulation, and allows
34 greater flexibility and ability to accommodate real datasets.

35 **Keywords:** area-interaction models; diagnostic tools; lasso; occupancy models;
36 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 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}_{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. The biases unique to the presence-only and presence-absence data sets are parameterised by $\boldsymbol{\alpha}_{\text{PO}}$ and $\boldsymbol{\alpha}_{\text{PA}}$, respectively, and collectively contained in the vector $\boldsymbol{\alpha}$. The key advancement of the combined likelihood approach is that the environmental response, parameterised by

68 β , is informed by both the presence-only and presence-absence data.

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

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

76 While these papers clearly advance the practice of fitting SDMs in important ways, they
77 do not address some common challenges that arise in real datasets. For example, they
78 all consider an inhomogeneous Poisson point process model (IPPPM) for the presence-
79 only data in the combination. In many real data sets, however, the implicit assumption
80 that the point locations are independently distributed conditional on the environment
81 is not met. Residual clustering or repulsion of the point locations not accounted for
82 with an IPPPM due to the observation process, unconsidered environmental covariates,
83 or biological factors would hence render the IPPPM inappropriate. Furthermore, none
84 of the current literature in combined likelihood approaches includes ways to account for
85 possible overfitting that results from including too many covariates in the model.

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

96 In this paper, we present a penalised combined likelihood model in a way that it is more
97 suitable for real data sets. In particular, we accommodate alternative forms of presence-
98 only models to account for spatial dependence and affix a penalty on model complexity
99 to address overfitting. In Section 2, we present the penalised combined likelihood formu-
100 lation. In Section 3, we illustrate via simulations the improvements that this formulation

101 provides and apply the proposed formulation to analyse the distribution of the Eurasian
 102 lynx (*lynx lynx*) in the Jura Mountains in eastern France. Finally, we present a discussion
 103 and further avenues for research in this area in Section 4.

104 2 Materials and Methods

105 2.1 Combined Penalised Likelihood Formulation

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

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

107 Here, $\boldsymbol{\alpha} = (\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_D)^\top$ is a q -dimensional vector that collects coefficients for the vari-
 108 ables \mathbf{Z} used to model bias for each of the D components individually. The environmental
 109 response is measured by a set of variables \mathbf{X} and is parametrised by $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$,
 110 which is collectively informed by all D components. The species data for all D compo-
 111 nents is collected in a set \mathbf{y} , with each individual data source \mathbf{y}_i determining the form of
 112 the component likelihood $\ell_i(\boldsymbol{\alpha}_i, \boldsymbol{\beta}; \mathbf{y}_i)$.

113 While many possibilities for the likelihood terms $\ell_i(\boldsymbol{\alpha}_i, \boldsymbol{\beta}; \mathbf{y}_i)$ are possible, we will focus
 114 on likelihood expressions for a PPM and for an occupancy model. For an IPPPM, we
 115 typically model the intensity of points $\mu(s)$ over a given study region \mathcal{A} as a log-linear
 116 function of environmental variables \mathbf{X} and bias terms \mathbf{Z} and derive estimates $\hat{\boldsymbol{\beta}}$ and $\hat{\boldsymbol{\alpha}}_{\text{PO}}$
 117 of the associated parameters by maximising a log-likelihood expression given by (Cressie,
 118 1992):

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

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

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

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

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

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

128 where A_i is the area of site i .

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

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

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

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

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

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

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

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

141 for $\gamma > 0$. Here, $\hat{\beta}_i^{(\text{unp})}$ is the unpenalised coefficient estimate corresponding to the i^{th}
142 environmental variable \mathbf{x}_i and $\hat{\alpha}_i^{(\text{unp})}$ is the unpenalised coefficient estimate corresponding
143 to the i^{th} bias variable \mathbf{z}_i . The shape of the weights is determined by the parameter γ .

144 The data-driven choice of the adaptive weights \mathbf{w} ensures that more important covariates
145 (*i.e.* those with coefficient estimates further away from 0) will be penalised less. This
146 construction also enables the adaptive lasso to achieve so-called oracle properties (Zou,
147 2006), which means that asymptotically, the correct subset of coefficients will be chosen
148 and the procedure has optimal estimation rate.

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

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

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

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

167 2.2 Implementation in R

168 To fit models with the combined penalised log-likelihood in eqn 1, we have developed
169 a set of functions in R inspired by the `optim` function and `ppmlasso` package (Renner
170 & Warton, 2013). The main function `comb.lasso` takes an input a list of species data,
171 associated environmental data, and formulae for the environmental trend and bias trends
172 for each component, along with details such as type of presence-only likelihoods to use,
173 the type of penalty, the number of models to fit, and the tuning parameter criterion.

174 The function applies the coordinate descent algorithm of Osborne *et al.* (2000). This
175 requires the derivatives of the component likelihoods (also known as “score equations”)
176 to be computed. Analytical score equations are supplied directly to the `optim` function,
177 which serves as the machinery of the optimisation. A tutorial illustrating use of this code
178 for the simulations as performed in Section 3.1 as well as some functions written to plot
179 intensity maps and features of the lasso penalisation is provided in the supplementary
180 material.

181 3 Results

182 3.1 Simulations

183 To investigate the benefits of the proposed penalised combined likelihood formulation, we
184 used the `rpoispp` function in `spatstat` (Baddeley & Turner, 2005) to generate a large
185 inhomogeneous Poisson pattern \mathbf{s}_{true} of roughly 10,000 points on a 30×30 -unit square
186 window from an intensity pattern defined by linear and quadratic terms of four generated
187 variables (hence eight meaningful covariates $\mathbf{x}_1, \dots, \mathbf{x}_8$).

188 From this pattern, we generated two biased presence-only subsamples \mathbf{s}_1 and \mathbf{s}_2 . The first
189 presence-only subsample \mathbf{s}_1 was biased by z_1 , the distance to a simulated road network,
190 and the other \mathbf{s}_2 by z_2 , the distance to a simulated categorical covariate. We varied the
191 size of the subsamples such that each pattern had either 50 or 200 points. We also varied
192 the strength of the clustering of the presence-only subsamples by setting the coefficient of
193 the interaction term β_{interact} . Here, the patterns either exhibit no clustering ($\beta_{\text{interact}} = 0$),
194 moderate clustering ($\beta_{\text{interact}} = 0.5$) or strong clustering ($\beta_{\text{interact}} = 1$). To sample the
195 points in \mathbf{s}_1 , we proceed as follows:

- 196 1. Initialise the set of sampled points $\mathbf{s}_1 = \emptyset$ and the point interactions to be a vector
197 of 0s
- 198 2. Compute the biased conditional intensity at every point in \mathbf{s}_{true} using $\mathbf{x}_1, \dots, \mathbf{x}_8$,
199 the bias covariate z_1 , and the current vector of point interactions
- 200 3. Set the conditional intensity for any point already selected in \mathbf{s}_1 to 0 to ensure these
201 points are not resampled
- 202 4. Randomly select a point from \mathbf{s}_{true} with sampling probabilities proportional to the
203 conditional intensities and add the selected point to \mathbf{s}_1

204 5. Update the vector of point interactions for all points in \mathbf{s}_{true} using the `evalInteraction`
205 function in `spatstat`

206 6. Repeat steps 2-5 until we have sampled the desired number of points

207 We sample \mathbf{s}_2 in a similar manner, using z_2 instead of z_1 to create the bias.

208 Because the true pattern \mathbf{s}_{true} is Poisson, this simulation setup emulates a scenario in which
209 the clustering of the observed point patterns is an artefact of the observation process –
210 this can happen if, for example, records are publicly available and enthusiasts for the
211 species report further observations near the publicly available locations (Johnston *et al.*,
212 2019).

213 We also generated a history \mathbf{y}_{occ} of detections and non-detections from 5 visits to each
214 of 100 sites centred along a regular grid in the 30×30 -unit observation window to
215 emulate a data set for which we could consider occupancy modelling. The species was
216 considered present at a site if the closest point in the pattern \mathbf{s}_{true} was within a distance
217 of 0.25 units of the centre of the site. The history of detections and non-detections at
218 each site where the species was considered present was randomly generated according
219 to detection probabilities defined by the inverse of the cloglog function evaluated at a
220 generated detection covariate \mathbf{z}_3 .

221 Finally, we generated 8 dummy covariates $\mathbf{d}_1, \dots, \mathbf{d}_8$ to include in fitted models that were
222 meaningless in describing the true species distribution. We did this to reflect the fact that
223 in real applications, we may not know which among a suite of candidate variables truly
224 determine the species distribution. We ensured that the maximum absolute correlation
225 among all pairs of variables was smaller than 0.5.

226 After generating the species data, we fit a number of models, using as input environmental
227 covariates the 8 meaningful covariates $\mathbf{x}_1, \dots, \mathbf{x}_8$ as well as 8 dummy covariates $\mathbf{d}_1, \dots, \mathbf{d}_8$
228 and using as bias covariates $\mathbf{z}_1, \mathbf{z}_2$, and \mathbf{z}_3 . For each of seven different combinations of
229 input data and presence-only likelihood, we fit a model without any penalty, with a lasso
230 penalty, and with an adaptive lasso penalty. For the models fitted with either a lasso
231 or an adaptive lasso penalty, we fit regularisation paths of 1000 models, increasing the
232 penalty from 0 to the smallest penalty λ_{max} that would shrink all coefficients to 0, thus
233 covering the entire scope of possible model sizes. The model which minimised BIC was
234 chosen among the 1000 fitted models. We considered as species data using a combination
235 of all three of $\mathbf{s}_1, \mathbf{s}_2$, and \mathbf{y}_{occ} as well as $\mathbf{s}_1, \mathbf{s}_2$, and \mathbf{y}_{occ} individually. For the models with
236 presence-only species inputs, we also varied whether these were modelled using an IPPPM

237 or an area-interaction model. This led to a total of 21 models being fitted, summarised
 238 in Table 1.

Model	Species Data	Presence-only likelihood	Penalty
1	$\mathbf{s}_1, \mathbf{s}_2,$ and \mathbf{y}_{occ}	IPPPM	None
2	$\mathbf{s}_1, \mathbf{s}_2,$ and \mathbf{y}_{occ}	IPPPM	Lasso
3	$\mathbf{s}_1, \mathbf{s}_2,$ and \mathbf{y}_{occ}	IPPPM	Adaptive Lasso
4	$\mathbf{s}_1, \mathbf{s}_2,$ and \mathbf{y}_{occ}	Area-interaction	None
5	$\mathbf{s}_1, \mathbf{s}_2,$ and \mathbf{y}_{occ}	Area-interaction	Lasso
6	$\mathbf{s}_1, \mathbf{s}_2,$ and \mathbf{y}_{occ}	Area-interaction	Adaptive Lasso
7	\mathbf{s}_1 only	IPPPM	None
8	\mathbf{s}_1 only	IPPPM	Lasso
9	\mathbf{s}_1 only	IPPPM	Adaptive Lasso
10	\mathbf{s}_1 only	Area-interaction	None
11	\mathbf{s}_1 only	Area-interaction	Lasso
12	\mathbf{s}_1 only	Area-interaction	Adaptive Lasso
13	\mathbf{s}_2 only	IPPPM	None
14	\mathbf{s}_2 only	IPPPM	Lasso
15	\mathbf{s}_2 only	IPPPM	Adaptive Lasso
16	\mathbf{s}_2 only	Area-interaction	None
17	\mathbf{s}_2 only	Area-interaction	Lasso
18	\mathbf{s}_2 only	Area-interaction	Adaptive Lasso
19	\mathbf{y}_{occ} only	–	None
20	\mathbf{y}_{occ} only	–	Lasso
21	\mathbf{y}_{occ} only	–	Adaptive Lasso

Table 1: Models fitted in each simulation. Models 1-6 were fitted using the proposed penalised combination framework, whereas models 7-21 were fitted using only one source of data. The models also varied based on the likelihood expression for any presence-only components and the type of penalty used, if any.

239 To evaluate performance, we compared the integrated mean squared error of the true in-
 240 tensity surface with rescaled fitted intensity surfaces of the 21 models. The fitted intensity
 241 surfaces were rescaled to have the same mean intensity as the true intensity surface to
 242 ensure that fair comparisons are made as models using different species data sources will
 243 have varying intercepts to reflect the estimated abundance of the points.

244 We performed 1,000 simulations of the data sets for each of the six combinations of

245 presence-only data set size and clustering strength and the resultant model fits on 512GB
246 nodes powered by 3.0 GHz Intel Xeon Gold (E5-6154) processor from the University
247 of Newcastle’s High Performance Computing cluster. The 6,000 simulation tasks took
248 approximately 12,000 hours.

249 Figure 1 shows boxplots of the calculated integrated mean squared errors from the simu-
250 lations. From these results, we can draw the following conclusions. First, the combined
251 likelihood approaches tend to outperform the models that rely on a single source of data,
252 confirming previous results by Dorazio (2014) and Fithian *et al.* (2015). This effect is
253 mitigated somewhat in the presence of clustering for the combined model which uses the
254 incorrect IPPPM likelihood, however.

255 Second, penalisation via the lasso or adaptive lasso improves model performance, and
256 this benefit is stronger with smaller data sets. This can be seen by comparing the same
257 models with different sample sizes (e.g. Models 7-9 when $N = 50$ to Models 7-9 when
258 $N = 200$ regardless of the clustering), as well as noting that the benefits of penalisation
259 are less pronounced for Models 1-3 and 4-6 which use all three data sets. This is an
260 expected conclusion given the danger of overfitting is greater with fewer observations.
261 Models penalised with the adaptive lasso tend to outperform models penalised with the
262 lasso for the presence-only data sets. Although the benefits of penalisation are negligible
263 with large data sets, fitting models with a penalty does not hurt the performance.

264 Finally, the performance benefits of selecting the correct presence-only area-interaction
265 likelihood in the presence of clustering tend to be greater with larger data sets. In the
266 first column of Figure 1, there is no clustering, such that models that use the IPPPM
267 likelihood for the presence-only components are appropriate, whereas in the second and
268 third columns, the models which use the area-interaction likelihood for the presence-only
269 components are appropriate. From the middle and right columns, we see that differences in
270 performance with the correct area-interaction likelihood for the presence-only components
271 are more pronounced as the presence-only sample size increases and as the strength of the
272 clustering becomes more pronounced. Curiously, the left column shows that the benefits
273 of the correct IPPPM likelihood are less pronounced with greater presence-only sample
274 sizes for Models 7-18 which only use a presence-only data source, and the combined models
275 appear to perform better with an area-interaction term when $N = 200$. We suspect that
276 the area-interaction term is not harming the performance because it is partially explained
277 by the effect of the bias used to generate the presence-only patterns – that is, biasing the
278 observations to be near roads or near level 1 of the categorical covariate induces some

279 clustering.

280 In summary, it appears that the proposed combined penalised likelihood framework pro-
281 vides the best performance.

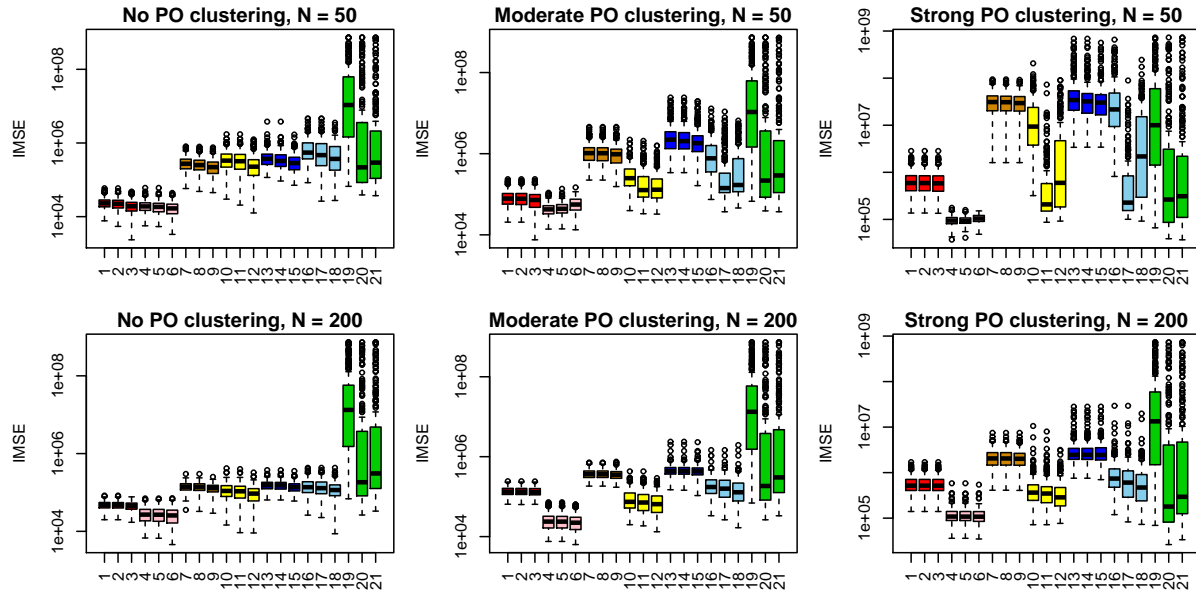


Figure 1: Boxplots of integrated mean squared error for the 21 models described in Table 1.

282 3.2 Analysis of Eurasian lynx distribution in the Jura Moun- 283 tains

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

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

293 We have three sources of lynx data in the Jura Mountains: a presence-only data set
294 consisting of 440 opportunistic sightings in the wild from 2009-2011 (denoted s_w), another
295 presence-only data set consisting of 240 reported interferences of lynx with domestic
296 livestock in 2009-2011 (denoted s_d), and pictures of lynx taken from cameras set up in

297 73 locations \mathbf{s}_c in the Jura Mountains in 2012. Lynx presence-only data were made
298 of presence signs sampled all year long thanks to a network of professional and non-
299 professional observers. Every observer is trained during a 3-day teaching course led by the
300 French National Game and Wildlife Agency (ONCFS) to document signs of the species'
301 presence (Duchamp *et al.*, 2012). Presence signs went through a standardised control
302 process to prevent misidentification (Duchamp *et al.*, 2012). The camera data has daily
303 reportings of the lynx across a total of 77 days. Due to this, we can consider the picture
304 history of lynx at the camera locations in an occupancy modelling framework (Blanc *et al.*,
305 2014). In particular, we split the 77-day period into seven 11-day periods, such that the
306 site history \mathbf{y}_c comprises seven detections and non-detections at each site in \mathbf{s}_c over each
307 11-day period.

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

314 To model the lynx distribution, we consider altitude, percentage of forest cover, distance
315 to the nearest water source, and human population density as environmental variables.
316 We model sampling bias in the wild records \mathbf{s}_w with distance to the nearest main road
317 and distance to the nearest train line, and sampling bias in the domestic records \mathbf{s}_d
318 with distance to the nearest farm and percentage of agricultural land. Finally, we model
319 detection probability for the camera data with distance to the nearest urban area. We
320 established this set of potential candidate environmental and detection variables based on
321 previously studied species habitat preferences and detectability (Bouyer *et al.*, 2015). The
322 Corine Land Cover land use repository from 2012 (Büttner *et al.*, 2014) supplies a map of
323 land coverage including urban areas, water areas, forest areas, farm areas, and agricultural
324 areas that was used to generate the percentage of forest areas and agricultural areas over
325 $1 \text{ km} \times 1 \text{ km}$ cells as well as distances to the nearest urban area, water source, and farm.
326 Altitude was averaged over $1 \text{ km} \times 1 \text{ km}$ cells from data available in the `raster` package
327 in R, while human population density was averaged over $1 \text{ km} \times 1 \text{ km}$ cells taken from
328 version 4 of the Gridded Population of the World data repository (Center for International
329 Earth Science Information Network (CIESIN) – Columbia University, 2016). Distances
330 from the nearest main road and railway were computed from shapefiles from Version 151

Lynx Locations

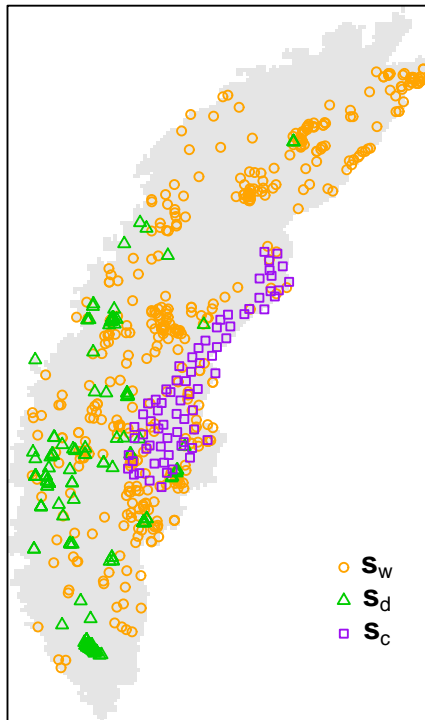


Figure 2: Locations of the lynx data in the Jura Mountains, including 440 observations in the wild s_w , 220 reports of domestic interference s_d , and 73 camera traps s_c .

331 of the ROUTE 500 database, accessible at <http://professionnels.ign.fr/route500>.

332 We fitted initial separate IPPPMs to the wild records s_w and the domestic records s_d
333 using linear, quadratic, and interaction terms for the four environmental covariates, and
334 linear terms for the bias covariates. From these models, we are able to assess whether
335 the assumption of independence inherent to the IPPPMs is appropriate with simulation
336 envelopes of the inhomogeneous K -function in `spatstat`, as shown in Figure 3. Both
337 of the envelopes for the IPPPMs fitted to the wild model (left panel) and the domestic
338 model (middle panel) demonstrate additional clustering as the observed inhomogeneous
339 K -function values plotted in red fall above the simulation envelopes for small radii. This
340 suggests that fitting an IPPPM is inappropriate for these data sets. The right panel
341 shows a simulation envelope comparing clustering across the two data sources, where the
342 intensities are estimated from area-interaction models, and as the observed values of the
343 K -function fall within the envelope boundaries, this suggests that there is no clustering
344 across the two data sets. This, in turn, suggests that the observed clustering within the
345 wild and domestic data sets may be more likely attributable to the observation process
346 than to some biological reality that induces clustering or a missed environmental covariate.

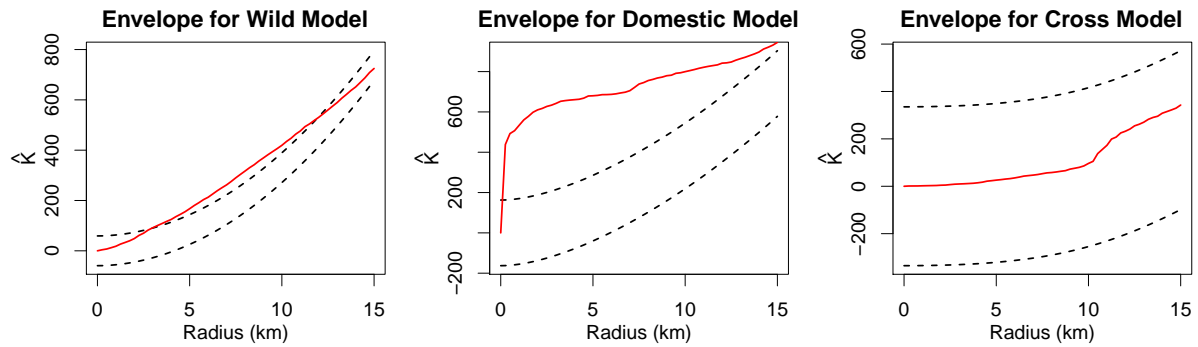


Figure 3: 95% simulation envelopes of the inhomogeneous K -function for the fitted IPPPM of the wild records (left), the fitted IPPPM of the domestic records (middle), and across the two fitted IPPPMs (right).

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

353 Figure 4 shows the bias-corrected fitted intensities from these two models. For the com-
354 bined model which uses IPPPMs (left panel), the fitted intensity is corrected for the
355 bias terms modelled for the presence-only components using the method of Warton *et al.*
356 (2013). For the combined penalised model which uses area-interaction models (right
357 panel), the fitted intensity is corrected for these same bias terms as well as the fitted
358 point interactions – that is, we treat the interaction parameter η as belonging to the set
359 of bias parameters α . The fitted models show strikingly different patterns, with the model
360 which uses area-interaction components highlighting much more of the Jura Mountains
361 as preferred habitat of lynx than the model which uses IPPPMs. We do not have access
362 to additional data with which to validate the performance of these models such as GPS
363 data as in Gould *et al.* (2019), but the results of Section 3.1 suggest that the model which
364 uses area-interaction components is likely to better reflect the true distribution of lynx.

365 The combined penalised model with the area-interaction components found the optimal
366 lasso penalty was 0, resulting in a model which included all 18 covariates and both inter-
367 action terms.

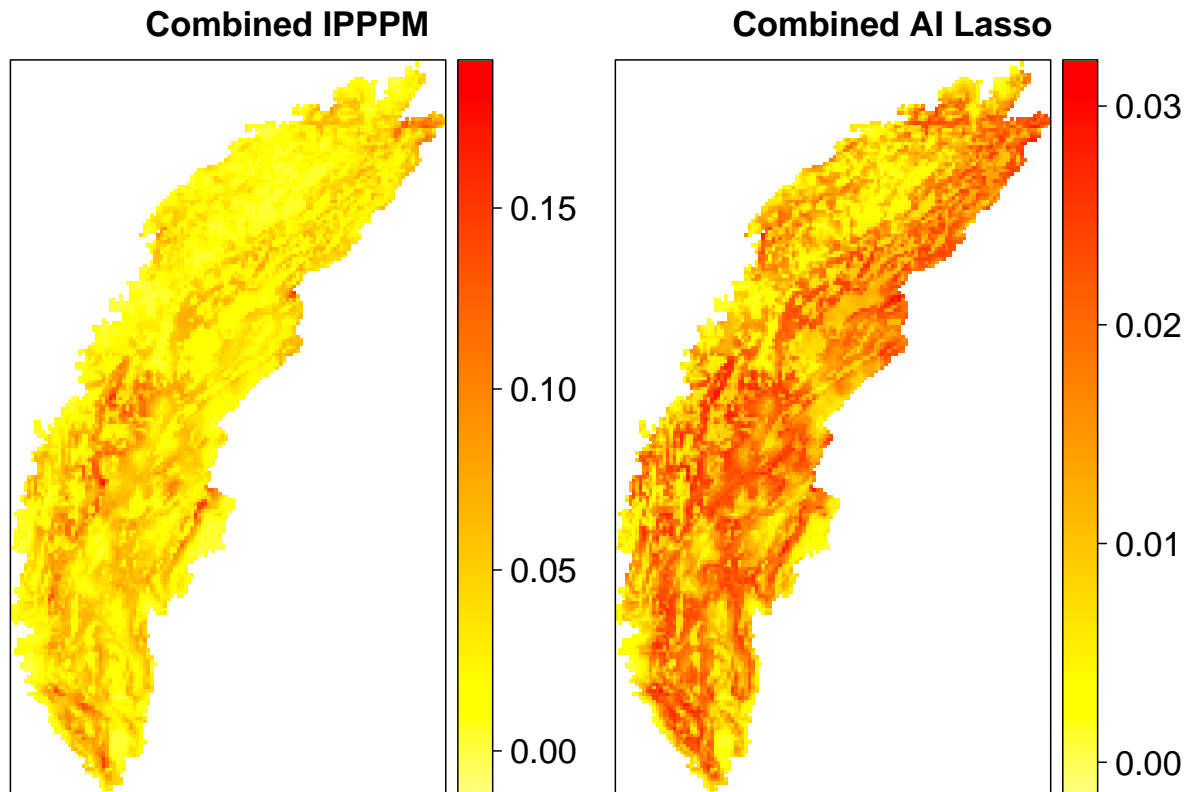


Figure 4: Fitted intensities using the combined likelihood formulation. Left: the model is fitted without any penalty and using inhomogeneous Poisson point process models for the presence-only data sources. Right: the model is fitted with a lasso penalty and using area-interaction models for the presence-only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

368 4 Discussion

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

373 4.1 Possible extensions

374 Despite these improvements, further advances are possible. Other penalty structures
375 could be incorporated into the same framework. While the lasso and adaptive lasso
376 showcased here show clear benefits in simulations, other penalised likelihood variants
377 such as SCAD (Fan & Li, 2001) could lead to superior performance in some situations,
378 and alternative methods to BIC of choosing the size of the penalty such as the Extended

379 Bayesian Information Criterion (“ERIC”, Hui *et al.*, 2015) could likewise be used.

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

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

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

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

399 4.2 Accounting for dependence within and among data sources

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

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

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

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

413 where $\xi(\mathbf{y}) \sim \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma})$. Models of this type are typically fitted in a Bayesian frame-
414 work. We could reduce the dimension of ξ through methods like fixed rank kriging or
415 induce sparsity in $\boldsymbol{\Sigma}$ through lasso-type penalties such that the likelihood in eqn 4 could
416 be fitted with software such as Template Model Builder (TMB, Kristensen *et al.*, 2016).

417 4.3 Conclusion and Perspectives

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

425 5 Acknowledgements

426 We thank the staff from the French National Game and Wildlife Agency, the Forest
427 National Agency and the Departmental Federation of Hunters of Jura department, who
428 collected the photographs during the camera-trapping session. We also thank all the
429 volunteers that are members of the Réseau Loup-Lynx that collect every year precious
430 presence signs of lynx all over its distribution area. We thank the University of Newcastle
431 through the Early to Mid-Career Researcher Visiting Fellowship and the Université Paul-
432 Valéry Montpellier through the Professeurs en mobilité universitaire fund for funding
433 visits for I.R. which helped facilitate the collaborations that resulted in this work.

434 6 Authors' Contributions

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

439 References

- 440 Baddeley, A. & Turner, R. (2000) Practical maximum pseudolikelihood for spatial point
441 patterns: (with discussion). *Australian & New Zealand Journal of Statistics*, **42**, 283–
442 322.
- 443 Baddeley, A. & Turner, R. (2005) Spatstat: an R package for analyzing spatial point
444 patterns. *Journal of Statistical Software*, **12**, 1–42.
- 445 Baddeley, A.J. & van Lieshout, M.N.M. (1995) Area-interaction point processes. *Annals*
446 *of the Institute of Statistical Mathematics*, **47**, 601–619.
- 447 Besag, J. (1977) Some methods of statistical analysis for spatial data. *Bulletin of the*
448 *International Statistical Institute*, **47**, 77–92.
- 449 Blanc, L., Marboutin, E., Gatti, S., Zimmermann, F. & Gimenez, O. (2014) Improving
450 abundance estimation by combining capture–recapture and occupancy data: example
451 with a large carnivore. *Journal of Applied Ecology*, **51**, 1733–1739.
- 452 Bouyer, Y., San Martin, G., Poncin, P., Beudels-Jamar, R.C., Odden, J. & Linnell, J.D.
453 (2015) Eurasian lynx habitat selection in human-modified landscape in norway: Effects
454 of different human habitat modifications and behavioral states. *Biological Conservation*,
455 **191**, 291–299.
- 456 Breitenmoser, U., Breitenmoser-Würsten, C. & Capt, S. (1998) Re-introduction and
457 present status of the lynx (*lynx lynx*) in switzerland. *Hystrix, the Italian Journal*
458 *of Mammalogy*, **10**.
- 459 Büttner, G., Soukup, T. & Kosztra, B. (2014) Clc2012 addendum to clc2006 technical
460 guidelines. *Final Draft, Copenhagen (EEA)*.
- 461 Center for International Earth Science Information Network (CIESIN) – Columbia Uni-
462 versity (2016) Gridded population of the world, version 4 (gpwv4): population density.

- 463 Cressie, N. (1992) *Statistics for spatial data*, volume 4. Wiley Online Library.
- 464 Cressie, N. & Wikle, C.K. (2015) *Statistics for spatio-temporal data*. John Wiley & Sons.
- 465 Diggle, P.J. (2003) *Statistical analysis of spatial point patterns*. Edward Arnold.
- 466 Dorazio, R.M. (2014) Accounting for imperfect detection and survey bias in statistical
467 analysis of presence-only data. *Global Ecology and Biogeography*, **23**, 1472–1484.
- 468 Duchamp, C., Boyer, J., Briaudet, P.E., Leonard, Y., Moris, P., Bataille, A., Dahier, T.,
469 Delacour, G., Millisher, G., Miquel, C. *et al.* (2012) A dual frame survey to assess time-
470 and space-related changes of the colonizing wolf population in france. *Hystrix-Italian*
471 *Journal of Mammalogy*, **23**, 14–28.
- 472 Fan, J. & Li, R. (2001) Variable selection via nonconcave penalized likelihood and its
473 oracle properties. *Journal of the American statistical Association*, **96**, 1348–1360.
- 474 Fithian, W., Elith, J., Hastie, T. & Keith, D.A. (2015) Bias correction in species dis-
475 tribution models: pooling survey and collection data for multiple species. *Methods in*
476 *Ecology and Evolution*, **6**, 424–438.
- 477 Gould, M.J., Gould, W.R., Cain III, J.W. & Roemer, G.W. (2019) Validating the perfor-
478 mance of occupancy models for estimating habitat use and predicting the distribution
479 of highly-mobile species: A case study using the american black bear. *Biological Con-*
480 *servation*, **234**, 28–36.
- 481 Guisan, A., Thuiller, W. & Zimmermann, N.E. (2017) *Habitat suitability and distribution*
482 *models: with applications in R*. Cambridge University Press.
- 483 Hirzel, A.H., Hausser, J., Chessel, D. & Perrin, N. (2002) Ecological-niche factor analysis:
484 how to compute habitat-suitability maps without absence data? *Ecology*, **83**, 2027–
485 2036.
- 486 Hui, F.K., Warton, D.I. & Foster, S.D. (2015) Tuning parameter selection for the adaptive
487 lasso using ERIC. *Journal of the American Statistical Association*, **110**, 262–269.
- 488 Hutchinson, R.A., Valente, J.J., Emerson, S.C., Betts, M.G. & Dietterich, T.G. (2015) Pe-
489 nalized likelihood methods improve parameter estimates in occupancy models. *Methods*
490 *in Ecology and Evolution*, **6**, 949–959.

- 491 Johnston, A., Hochachka, W., Strimas-Mackey, M., Gutierrez, V.R., Robinson, O., Miller,
492 E., Auer, T., Kelling, S. & Fink, D. (2019) Best practices for making reliable inferences
493 from citizen science data: case study using ebird to estimate species distributions.
494 *bioRxiv*, p. 574392.
- 495 Koshkina, V., Wang, Y., Gordon, A., Dorazio, R.M., White, M. & Stone, L. (2017)
496 Integrated species distribution models: combining presence-background data and site-
497 occupancy data with imperfect detection. *Methods in Ecology and Evolution*, **8**, 420–430.
- 498 Kristensen, K., Nielsen, A., Berg, C.W., Skaug, H. & Bell, B.M. (2016) TMB: Automatic
499 differentiation and Laplace approximation. *Journal of Statistical Software*, **70**, 1–21.
- 500 MacKenzie, D.I., Nichols, J.D., Hines, J.E., Knutson, M.G. & Franklin, A.B. (2003)
501 Estimating site occupancy, colonization, and local extinction when a species is detected
502 imperfectly. *Ecology*, **84**, 2200–2207.
- 503 Miller, D.A., Pacifici, K., Sanderlin, J.S. & Reich, B.J. (2019) The recent past and promis-
504 ing future for data integration methods to estimate species' distributions. *Methods in*
505 *Ecology and Evolution*, **10**, 22–37.
- 506 Osborne, M.R., Presnell, B. & Turlach, B.A. (2000) On the lasso and its dual. *Journal*
507 *of Computational and Graphical Statistics*, **9**, 319–337.
- 508 Pacifici, K., Reich, B.J., Miller, D.A., Gardner, B., Stauffer, G., Singh, S., McKerrow,
509 A. & Collazo, J.A. (2017) Integrating multiple data sources in species distribution
510 modeling: A framework for data fusion. *Ecology*, **98**, 840–850.
- 511 Papangelou, F. (1974) The conditional intensity of general point processes and an appli-
512 cation to line processes. *Probability Theory and Related Fields*, **28**, 207–226.
- 513 Phillips, S.J., Anderson, R.P. & Schapire, R.E. (2006) Maximum entropy modeling of
514 species geographic distributions. *Ecological Modelling*, **190**, 231–259.
- 515 Renner, I.W. & Warton, D.I. (2013) Equivalence of MAXENT and Poisson point process
516 models for species distribution modeling in ecology. *Biometrics*, **69**, 274–281.
- 517 Renner, I.W., Elith, J., Baddeley, A., Fithian, W., Hastie, T., Phillips, S.J., Popovic,
518 G. & Warton, D.I. (2015) Point process models for presence-only analysis. *Methods in*
519 *Ecology and Evolution*, **6**, 366–379.

- 520 Tibshirani, R. (1996) Regression shrinkage and selection via the lasso. *Journal of the*
521 *Royal Statistical Society Series B (Methodological)*, pp. 267–288.
- 522 Vandiel, J.M. & Stahl, P. (2005) Distribution trend of the eurasian lynx lynx lynx popu-
523 lations in france. *Mammalia mamm*, **69**, 145–158.
- 524 Warton, D.I., Renner, I.W. & Ramp, D. (2013) Model-based control of observer bias for
525 the analysis of presence-only data in ecology. *PloS one*, **8**, e79168.
- 526 Warton, D.I. & Shepherd, L.C. (2010) Poisson point process models solve the “pseudo-
527 absence problem” for presence-only data in ecology. *Annals of Applied Statistics*, **4**,
528 1383–1402.
- 529 Widom, B. & Rowlinson, J.S. (1970) New model for the study of liquid–vapor phase
530 transitions. *The Journal of Chemical Physics*, **52**, 1670–1684.
- 531 Zou, H. (2006) The adaptive lasso and its oracle properties. *Journal of the American*
532 *Statistical Association*, **101**, 1418–1429.