

Spatial occupancy models for data collected on stream networks

Olivier Gimenez^{1*}

2024-09-09

¹ CEFE, Univ Montpellier, CNRS, EPHE, IRD, Montpellier, France

* Corresponding author: olivier.gimenez@cefe.cnrs.fr

To effectively monitor biodiversity in streams and rivers, we need to quantify species distribution accurately. Occupancy models are useful for distinguishing between the non-detection of a species and its actual absence. While these models can account for spatial autocorrelation, they are not suited for streams and rivers due to their unique network spatial structure. Here, I propose spatial occupancy models specifically designed for data collected on stream and river networks. I present the statistical developments and illustrate their application using data on a semi-aquatic mammal. Overall, spatial stream network occupancy models offer a robust method for assessing biodiversity in freshwater ecosystems.

Keywords: Bayesian statistics, Spatial stream network models, Occupancy models, Spatial autocorrelation, Wildlife monitoring

18 Introduction

19 Streams and rivers provide essential habitats for numerous species of animals and plants ([Reid](#)
20 [et al. 2019](#)). The ecological health of these freshwater ecosystems is paramount not only for the
21 biodiversity they harbor but also for the ecosystem services they provide, which are
22 indispensable to both wildlife and human populations ([Vári et al. 2021](#)). However, human
23 activities are altering the natural conditions of streams, rivers, and their associated riparian
24 habitats, jeopardizing the persistence of these ecosystems ([Albert et al. 2020](#)).

25 In this context, species distribution models (SDMs) are essential tools in understanding and
26 preserving biodiversity ([Elith and Leathwick 2009](#)). SDMs predict the distribution of species,
27 helping scientists and conservationists in identifying critical habitats and biodiversity hotspots.
28 Additionally, SDMs inform strategies aimed at mitigating the impacts of climate and land-use
29 changes, managing invasive species, and enhancing habitat connectivity in freshwater
30 ecosystems ([Domisch et al. 2015](#)).

31 SDMs are influenced by two main issues: imperfect detection and spatial autocorrelation
32 ([Guélat and Kéry 2018](#)). First, imperfect detection occurs when a species present in a given area
33 is not detected during surveys due to factors such as observer experience, species behavior, and
34 environmental conditions. Ignoring imperfect detection can lead to biased estimates of species
35 distribution and flawed inferences about the relationship between species presence and
36 environmental factors (e.g., [Lahoz-Monfort et al. 2014](#)). This can misinform conservation
37 strategies and habitat management decisions. To address this issue, occupancy models are
38 SDMs that rely on repeated visits of spatial sampling units for inferring distribution
39 ([MacKenzie et al. 2017](#)). These models have been widely used in freshwater ecosystems for
40 various taxa (e.g., [Wedderburn et al. 2022](#), [Couturier et al. 2023](#)).

41 Second, SDMs rely on the assumption of independent residuals. This assumption may be
42 violated if nearby sampling sites tend to have similar probabilities of species presence, leading

to biased estimates of species distribution and potentially inflating the effects of environmental factors (Dormann et al. 2007). Several extensions of occupancy models have been proposed to account for spatial autocorrelation, building on the spatial statistics literature, with conditional autoregressive models (Johnson et al. 2013) and geospatial models (Rushing et al. 2019). However, these models rely on the Euclidean distance between the spatial sampling units, which does not acknowledge the spatial structure in networks of streams and rivers.

Here I propose spatial occupancy models that account for spatial autocorrelation based on flow connectivity and river network (Peterson et al. 2013). I build on the linear mixed modelling approach developed by Ver Hoef and Peterson (2010) and Peterson and Hoef (2010), which integrates various distance-based spatial correlation structures (both Euclidean and non-Euclidean) within a single model. I plug-in this variance component approach into occupancy models using a Bayesian approach. A similar approach was recently undertaken by Lu et al. (2024) for count data to estimate abundance. Below I outline the statistical developments of this model and demonstrate its application to a semi-aquatic mammal in French streams and rivers.

Methods

Occupancy models

To address imperfect detection, I use occupancy models to estimate the true species distribution (MacKenzie et al. 2017). In these models, monitoring occurs across S spatial sampling units, or sites. If detection was perfect, the state z_i of a site i would be a Bernoulli random variable, taking value 1 with occupancy probability ψ_i if the site was occupied, and 0 otherwise with probability $1 - \psi_i$. However, because the ecological process of state occupancy z_i is only partially observable (since the species might be present but undetected), we must also account

for the observation process, which is also modeled as a Bernoulli random variable. When the species is detected at site i , i.e. $y_i = 1$, with detection probability p , it confirms that the site is occupied. Conversely, if the species is not detected, i.e. $y_i = 0$, with probability $1 - p$, we cannot determine whether the site is occupied or not. Both parameters, ψ and p , can be modeled as functions of explanatory spatial variables, in the spirit of generalized linear models, and logistic regressions for example.

Spatial autocorrelation for stream networks

How is spatial autocorrelation accounted for in occupancy models? The usual way is to write the probabilities of occupancy $\psi = (\psi_1, \dots, \psi_S)$ on some scale, say the logit scale, as a function of P explanatory variables gathered in a matrix \mathbf{X} with S rows and $P + 1$ columns with a corresponding vector of regression parameters $\beta = (\beta_0, \beta_1, \dots, \beta_P)$ that need to be estimated. To account for spatial autocorrelation, a random effect ϵ is added to the model, which captures the spatial dependencies among sites (Guélat and Kéry 2018):

$$\text{logit}(\psi) = \mathbf{X}\beta + \epsilon. \quad (1)$$

The random effect ϵ can be structured using methods such as conditional autoregressive models and their extensions (Johnson et al. 2013), or geoaddivitive models (Rushing et al. 2019). However, these approaches typically rely on Euclidean distance to assess proximity among sites, which may not fully capture the complex spatial dependencies present in streams and rivers. Specifically, we are interested in flow connectivity and the topology of streams and rivers (Peterson et al. 2013). Following Ver Hoef and Peterson (2010) and Peterson and Hoef (2010), I define two sites as flow-connected if water flows from an upstream site to a downstream site, and as flow-unconnected if they share a common confluence downstream but do not directly share flow. Then I parameterize occupancy by rewriting the random effect as a

mixture of four components as follows:

$$\text{logit}(\psi) = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\tau}_{tu} + \boldsymbol{\tau}_{td} + \boldsymbol{\tau}_{eu} + \epsilon \quad (2)$$

where $\boldsymbol{\tau}_{tu}$ is a random effect with spatial covariance between flow-connected sites that can occur in the same direction of the river flow (tail-up), which is suitable for organisms that move passively, such as mussels), $\boldsymbol{\tau}_{td}$ is a random effect with spatial covariance between flow-connected and flow-unconnected sites, which can occur with or against the direction of flow (tail-down), applicable to organisms that move actively, such as semi-aquatic mammals, $\boldsymbol{\tau}_{eu}$ is a random effect with a spatial covariance independent of the network topology, influenced by factors like air temperature or precipitation, and ϵ is a random effect with variance, often referred to as the nugget, which accounts for additional variability. How to build these covariance components is described in details elsewhere (e.g., [Ver Hoef et al. 2019](#)), and I provide an example in the next section.

Case study

To illustrate the new approach, I investigated the impact of human disturbance on the occupancy of European otter (*Lutra lutra*) in France. The otter, a semi-aquatic mammal, faced near extinction in the 20th century in France due to extensive hunting for its fur. With hunting bans and protection efforts, the species is now recolonizing the country, and the ecological question is assessing its current distribution. Data on otter detection and non-detection were collected in 2003-2005 in the Midi-Pyrénées region (see panel a in Fig. 1). Observers searched for signs of otter presence at a small river catchment scale, which was used as the spatial sampling unit. These data were analyzed by Couturier et al. (2023), who found that human density and the proportion of cultivated area influenced occupancy. In this study, I focus on a subsample of this dataset, covering $S = 56$ sites in the Lot, Aveyron and Cantal counties, which

110 were visited 3 times (see panel b in Fig. 1). I considered $P = 2$ explanatory variables. I used
 111 human population density as a proxy for human disturbance, calculated as the number of
 112 inhabitants per km² within a 200-m buffer around each stream (see panel c in Fig. 1).
 113 Additionally, I considered the proportion of cultivated areas. Detection was considered as
 114 constant. Regarding spatial autocorrelation, since otters can move both downstream and
 115 upstream, I used a tail-down model:

$$\text{logit}(\psi) = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\tau}_{td} \quad (3)$$

116 with an exponential covariance:

$$\text{Cov}(\tau_i, \tau_j) = \begin{cases} \sigma^2 \exp(-h/\theta), & \text{if sites } i \text{ and } j \text{ are flow-connected} \\ \sigma^2 \exp(-(a+b)/\theta), & \text{if sites } i \text{ and } j \text{ are flow-unconnected} \end{cases}$$

117 where I dropped the td notation in τ_{td} for clarity, h is the stream distance between sites i and j , b
 118 denotes the longer of the distances between sites i and j to the common downstream junction, a
 119 denotes the shorter of the two distances, σ^2 is a variance parameter usually referred to as the
 120 partial sill and θ is a range parameter (see Appendix A in [Peterson and Hoef 2010](#)).

121 **Simulation study**

122 I conducted a simulation study to evaluate model performance by examining parameter bias
 123 and prediction error. I simulated a stream network over $S = 100$ sites using an exponential
 124 tail-down structure (inspired by the case study, see previous section) with partial sill $\sigma^2 = 2$
 125 and range parameter $\theta = 10$. I considered a single covariate, normally distributed with mean 0
 126 and standard deviation 1, with a linear effect (on the logit scale) on the occupancy probability
 127 with intercept $\beta_0 = 0.5$ and slope $\beta_1 = 1$. I simulated the observation process with a detection

probability $p = 0.6$ across 5 repeated visits per site. I fitted both models, with and without spatial autocorrelation, to the simulated data, and I repeated this procedure 100 times. Eventually, I calculated the relative bias for each parameter and the root mean square prediction error (RMSPE) for each model.

Bayesian implementation

For all analyses, I used the statistical language R (R Core Team 2023). I used the tidyverse (Wickham et al. 2019) suite of packages for data manipulation and visualization and sf (Pebesma and Bivand 2023) for dealing with spatial data. I used the openSTARS (Kattwinkel et al. 2020) and SSN (Ver Hoef et al. 2014) packages to build and characterize the network and calculate stream distances. To implement the model, I was inspired by code for spatial stream network models provided in Santos-Fernandez et al. (2022) and code for occupancy models from Clark and Altwegg (2019). I fitted models within the Bayesian framework using the rstan (Stan Development Team 2024) package. I post-processed the results with the bayesplot (Gabry et al. 2019) package. I specified weakly informative priors for all parameters, specifically normal distributions with mean 0 and standard deviation 1.5 for regression parameters, and uniform distributions for the partial sill, range parameter and detection probability. I ran two chains for a total of 15,000 iterations with a burn-in of 5,000 iterations. I summarized posterior distributions with posterior mean and 95% credible intervals. I assessed model convergence using R-hat values (< 1.1), effective sample size (> 100), and visual inspection of the trace plots.

Results and discussion

Here, I provide the parameter estimates from the new model accommodating spatial autocorrelation (see also Fig. 2), unless otherwise specified. Detection probability was less than one, estimated at 0.71 (0.59, 0.80), which justified the use of occupancy models. The proportion

of cultivated area had no effect, with a slope estimated at 0.60 (-0.67, 1.96). Population density also had no effect on occupancy probability, with a slope estimated at -0.96 (-2.24, 0.17).

However, I did find a negative effect when spatial autocorrelation was ignored, with a slope estimated at -1.10 (-1.99, -0.34). This latter result aligns with a previous analysis of a more extensive dataset (Couturier et al. 2023) that also ignored spatial autocorrelation.

As anticipated, the effect size of human density increased when spatial autocorrelation was ignored. The most likely explanation for this is a bias due to an omitted variable. There is spatial autocorrelation in human density (see panel c in Fig. 1), which inflates its effect size; this bias is controlled for when spatial autocorrelation is included in the model. There must be spatial variation in occupancy probabilities attributable to another variable that needs to be accounted for.

The results of the simulation study underscored the importance of accounting for spatial autocorrelation. Ignoring spatial autocorrelation led to a relative bias of -26% in the slope of the covariate, compared to just 0.7% when spatial autocorrelation was included. Although the range parameter exhibited substantial bias (330%), this outcome was expected. In spatial models, there is typically a strong positive relationship between the range parameter and the partial sill (Zhang 2004, Ver Hoef et al. 2024), which showed a bias of 12%. Importantly, the ratio of the range parameter to the partial sill can be reliably estimated (Zhang 2004), as evidenced by its negligible bias (0.9%). Overall, the root mean square prediction error (RMSPE) was lower for the model with spatial correlation (RMSPE = 0.18) compared to the model without it (RMSPE = 0.23).

Two short-term perspectives arise from this work. From a methodological perspective, the new approach could be extended to multi-season occupancy models, enabling the modeling of colonization probability as a function of distance to habitat features that may impede species movement. This would facilitate the quantification of landscape connectivity in freshwater

ecosystems. Such development requires moving to spatio-temporal models for stream and river data, which have recently become available ([Santos-Fernandez et al. 2022](#)). From an ecological perspective, the new approach presents significant potential for the analysis of environmental DNA (eDNA). The eDNA methodology offers substantial promise for the non-invasive monitoring of biodiversity in freshwater ecosystems ([Carraro et al. 2020](#)). While spatial stream network models have been employed to analyze eDNA data ([Winkowski et al. 2024](#)), these models have overlooked the issue of imperfect detection. Previous studies have recognized occupancy models as effective tools for eDNA data analysis ([Burian et al. 2021](#)), with some considering spatial autocorrelation ([Chen and Ficetola 2019](#)), however they have yet to integrate spatial stream networks. The new approach addresses this gap by incorporating both imperfect detection and spatial stream networks.

Acknowledgments

First, I would like to warmly thank Jay Ver Hoef and Edgar Santos Fernández for useful discussions on spatial stream network models. Also, I thank Maëlis Kervellec for sharing her code, and Thibaut Couturier for explaining the otter data. Finally, I thank all structures and people who participated in the coordination, fieldwork and data management of the otter monitoring.

Ethics and Integrity statements

Not applicable.

Data availability statement

Data and code are available at

<https://github.com/oliviergimenez/spatial-stream-network-occupancy-model>.

Funding statement

This research is a product of the DISCAR group funded by the French Foundation for Research on Biodiversity (FRB) through its synthesis center CESAB.

Conflict of interest disclosure

The author has no conflicts of interest to declare.

References

- Albert, J. S., G. Destouni, S. M. Duke-Sylvester, A. E. Magurran, T. Oberdorff, R. E. Reis, K. O. Winemiller, and W. J. Ripple. 2020. Scientists' warning to humanity on the freshwater biodiversity crisis. *Ambio* 50:85–94.
- Burian, A., Q. Mauvisseau, M. Bulling, S. Domisch, S. Qian, and M. Sweet. 2021. Improving the reliability of eDNA data interpretation. *Molecular Ecology Resources* 21:1422–1433.
- Carraro, L., E. Mächler, R. Wüthrich, and F. Altermatt. 2020. Environmental DNA allows upscaling spatial patterns of biodiversity in freshwater ecosystems. *Nature Communications* 11:3585.
- Chen, W., and G. F. Ficetola. 2019. Conditionally autoregressive models improve occupancy analyses of autocorrelated data: An example with environmental DNA. *Molecular Ecology Resources* 19:163–175.

215 Clark, A. E., and R. Altwegg. 2019. Efficient Bayesian analysis of occupancy models with logit
216 link functions. *Ecology and Evolution* 9:756–768.

217 Couturier, T., J. Steinmetz, P. Defos du Rau, D. Marc, E. Trichet, R. Gomes, and A. Besnard. 2023.
218 Intensive agriculture as the main limiting factor of the otter's return in southwest france.
219 *Biological Conservation* 279:109927.

220 Domisch, S., S. Jähnig, J. Simaika, M. Kuemmerlen, and S. Stoll. 2015. Application of species
221 distribution models in stream ecosystems: The challenges of spatial and temporal scale,
222 environmental predictors and species occurrence data. *Fundamental and Applied*
223 *Limnology* 186:45–61.

224 Dormann, F. C., J. M. McPherson, M. B. Araújo, R. Bivand, J. Bolliger, G. Carl, R. G. Davies, A.
225 Hirzel, W. Jetz, W. Daniel Kissling, I. Kühn, R. Ohlemüller, P. R. Peres-Neto, B. Reineking, B.
226 Schröder, F. M. Schurr, and R. Wilson. 2007. Methods to account for spatial autocorrelation
227 in the analysis of species distributional data: A review. *Ecography* 30:609–628.

228 Elith, J., and J. R. Leathwick. 2009. Species distribution models: Ecological explanation and
229 prediction across space and time. *Annual Review of Ecology, Evolution, and Systematics*
230 40:677–697.

231 Gabry, J., D. Simpson, A. Vehtari, M. Betancourt, and A. Gelman. 2019. [Visualization in](#)
232 [bayesian workflow](#). *J. R. Stat. Soc. A* 182:389–402.

233 Guélat, J., and M. Kéry. 2018. Effects of spatial autocorrelation and imperfect detection on
234 species distribution models. *Methods in Ecology and Evolution* 9:1614–1625.

235 Johnson, D. S., P. B. Conn, M. B. Hooten, J. C. Ray, and B. A. Pond. 2013. Spatial occupancy
236 models for large data sets. *Ecology* 94:801–808.

237 Kattwinkel, M., E. Szöcs, E. Peterson, and R. Schäfer. 2020. Preparing GIS data for analysis of
238 stream monitoring data: The R package openSTARS. *Plos One* 15:e0239237.

239 Lahoz-Monfort, J. J., G. Guillera-Arroita, and B. A. Wintle. 2014. Imperfect detection impacts
240 the performance of species distribution models. *Global Ecology and Biogeography*

23:504–515.

Lu, X., Y. Kanno, G. P. Valentine, J. M. Rash, and M. B. Hooten. 2024. Using multi-scale spatial models of dendritic ecosystems to infer abundance of a stream salmonid. *Journal of Applied Ecology* 61:1703–1715.

MacKenzie, D. I., J. D. Nichols, J. A. Royle, K. H. Pollock, L. L. Bailey, and J. E. Hines. 2017. *Occupancy Estimation and Modeling: Inferring Patterns and Dynamics of Species Occurrence*. Elsevier.

Pebesma, E., and R. Bivand. 2023. [Spatial Data Science: With applications in R](#). Chapman and Hall/CRC.

Peterson, E. E., and J. M. V. Hoef. 2010. A mixed-model moving-average approach to geostatistical modeling in stream networks. *Ecology* 91:644–651.

Peterson, E. E., J. M. Ver Hoef, D. J. Isaak, J. A. Falke, M.-J. Fortin, C. E. Jordan, K. McNyset, P. Monestiez, A. S. Ruesch, A. Sengupta, N. Som, E. A. Steel, D. M. Theobald, C. E. Torgersen, and S. J. Wenger. 2013. Modelling dendritic ecological networks in space: An integrated network perspective. *Ecology Letters* 16:707–719.

R Core Team. 2023. [R: A language and environment for statistical computing](#). R Foundation for Statistical Computing, Vienna, Austria.

Reid, A. J., A. K. Carlson, I. F. Creed, E. J. Eliason, P. A. Gell, P. T. J. Johnson, K. A. Kidd, T. J. MacCormack, J. D. Olden, S. J. Ormerod, J. P. Smol, W. W. Taylor, K. Tockner, J. C. Vermaire, D. Dudgeon, and S. J. Cooke. 2019. Emerging threats and persistent conservation challenges for freshwater biodiversity. *Biological Reviews* 94:849–873.

Rushing, C., J. Andrew Royle, D. Ziolkowski Jr, and K. Pardieck. 2019. Modeling spatially and temporally complex range dynamics when detection is imperfect. *Scientific Reports* 9.

Santos-Fernandez, E., J. M. Ver Hoef, E. E. Peterson, J. McGree, D. J. Isaak, and K. Mengersen. 2022. Bayesian spatio-temporal models for stream networks. *Computational Statistics & Data Analysis* 170:107446.

267 Stan Development Team. 2024. [RStan: The R interface to Stan](#).

268 Vári, Á., S. Podschun, T. Eros, T. Hein, B. Pataki, C. Ioja, C. Adamescu, A. Gerhardt, T. Gruber,
 269 A. Dedić, M. Ciric, B. Gavrilović, and A. Báldi. 2021. Freshwater systems and ecosystem
 270 services: Challenges and chances for cross-fertilization of disciplines. *Ambio* 51:135–151.

271 Ver Hoef, J. M., E. Blagg, M. Dumelle, P. M. Dixon, D. L. Zimmerman, and P. B. Conn. 2024.
 272 Marginal inference for hierarchical generalized linear mixed models with patterned
 273 covariance matrices using the laplace approximation. *Environmetrics* n/a:e2872.

274 Ver Hoef, J., and E. Peterson. 2010. A moving average approach for spatial statistical models of
 275 stream networks. *Journal of the American Statistical Association* 105:6–18.

276 Ver Hoef, J., E. Peterson, D. Clifford, and R. Shah. 2014. SSN: An R package for spatial statistical
 277 modeling on stream networks. *Journal of Statistical Software* 56:1–45.

278 Ver Hoef, J., E. Peterson, and D. Isaak. 2019. Spatial statistical models for stream networks.
 279 Pages 421–441 in A. E. Gelfand, M. Fuentes, J. A. Hoeting, and S. R. Lyttleton, editors.
 280 Handbook of environmental and ecological statistics. Chapman; Hall/CRC.

281 Wedderburn, S. D., N. S. Whiterod, and L. Vilizzi. 2022. Occupancy modelling confirms the first
 282 extirpation of a freshwater fish from one of the world’s largest river systems. *Aquatic
 283 Conservation: Marine and Freshwater Ecosystems* 32:258–268.

284 Wickham, H., M. Averick, J. Bryan, W. Chang, L. D. McGowan, R. François, G. Grolemond, A.
 285 Hayes, L. Henry, J. Hester, M. Kuhn, T. L. Pedersen, E. Miller, S. M. Bache, K. Müller, J.
 286 Ooms, D. Robinson, D. P. Seidel, V. Spinu, K. Takahashi, D. Vaughan, C. Wilke, K. Woo, and
 287 H. Yutani. 2019. [Welcome to the tidyverse](#). *Journal of Open Source Software* 4:1686.

288 Winkowski, J. J., J. D. Olden, and S. Brown. 2024. Integrating spatial stream network models
 289 and environmental DNA to estimate current and future distributions of nonnative
 290 smallmouth bass. *Transactions of the American Fisheries Society* 153:180–199.

291 Zhang, H. 2004. Inconsistent estimation and asymptotically equal interpolations in
 292 model-based geostatistics. *Journal of the American Statistical Association* 99:250–261.

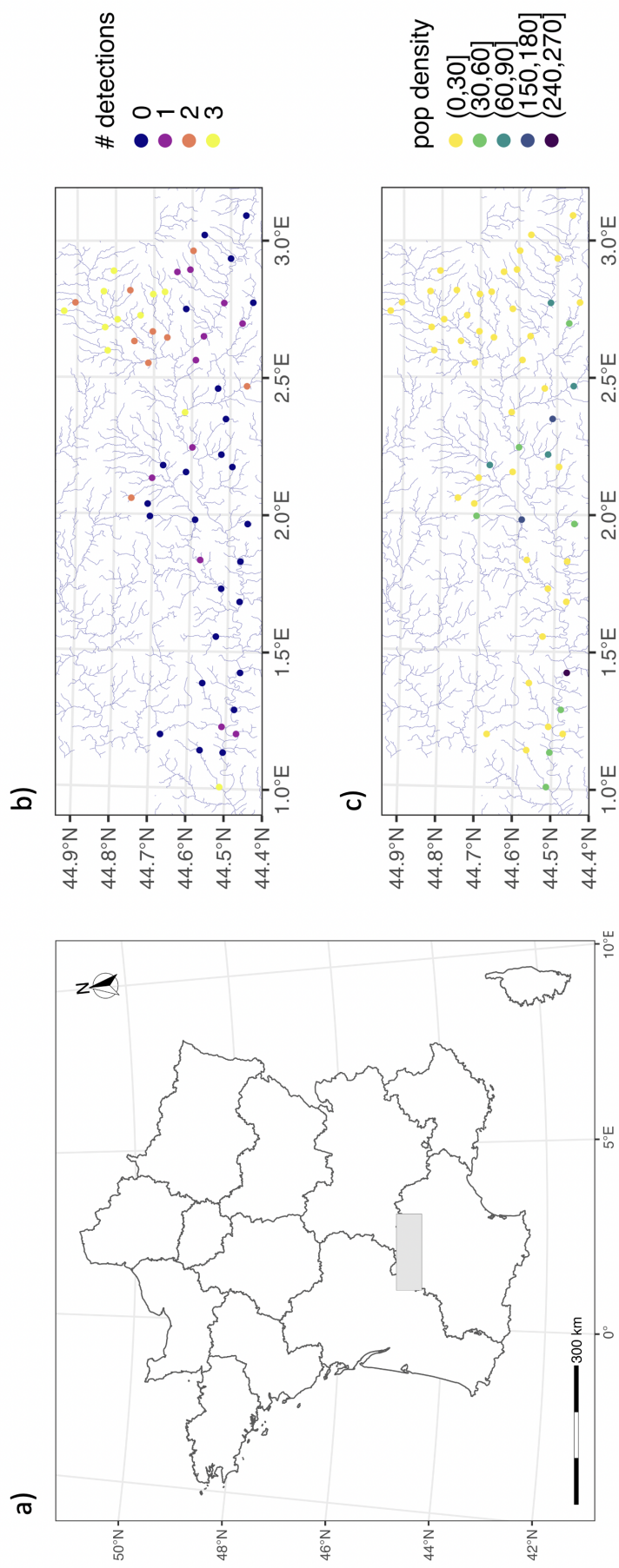


Figure 1: Information on the otter data is provided. In panel a), the study area is given in a grey rectable on a map of France. In panel b) the number of detections is given on a map of the study area; it is the sum of detections across the three repeated visits, which can range from 0 (if the species was not detected at any visit) to 3 (if the species was detected at all three visits). In panel c) the human population density is represented on a map of the study area, in number of inhabitants per km²; the missing intervals (90,120] and (120,150] are due to the absence of sites with human population densities within these ranges.

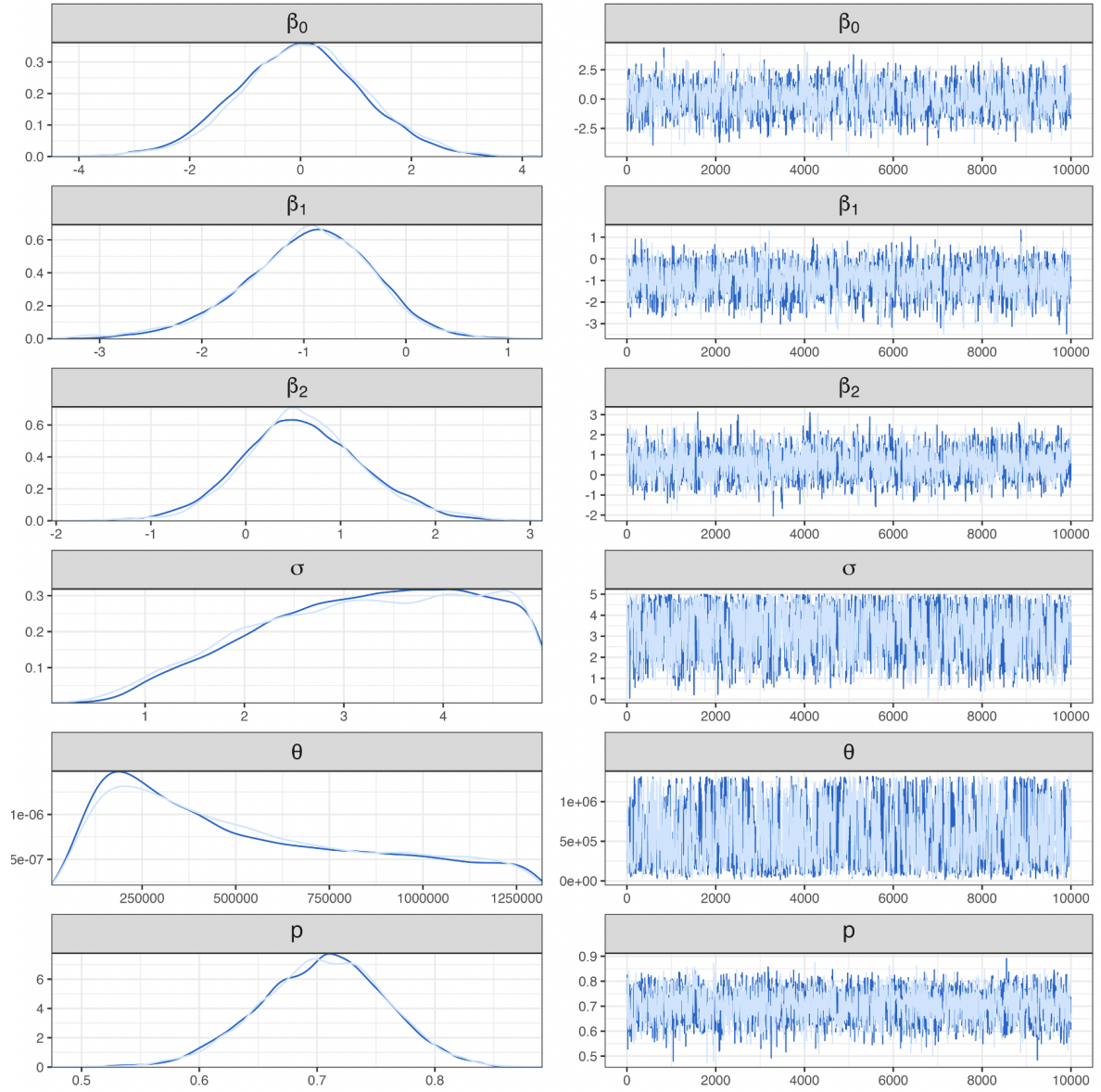


Figure 2: Posterior density plots (left) and trace plots (right) are shown for the regression parameters (β_0 : intercept, β_1 : slope of the human population density effect, β_2 : slope of the proportion of cultivated areas effect), spatial parameters (σ : partial sill, θ : range parameter) and the detection probability (p). Each of the two chains is represented in a distinct shade of blue.