

# Central Lancashire Online Knowledge (CLoK)

| Title    | Bayesian Estimation of Heterogeneous Environments from Animal            |
|----------|--|
|          | Movement Data  |
| Туре     | Article  |
| URL      | https://clok.uclan.ac.uk/id/eprint/37538/                                |
| DOI      | https://doi.org/10.1002/env.2679   |
| Date     | 2021   |
| Citation | Tishkovskaya, Svetlana and Blackwell, Paul G. (2021) Bayesian Estimation |
|          | of Heterogeneous Environments from Animal Movement Data.                 |
|          | Environmetrics, 32 (6). e2679. ISSN 1180-4009                            |
| Creators | Tishkovskaya, Svetlana and Blackwell, Paul G.                            |

It is advisable to refer to the publisher's version if you intend to cite from the work. https://doi.org/10.1002/env.2679

For information about Research at UCLan please go to <a href="http://www.uclan.ac.uk/research/">http://www.uclan.ac.uk/research/</a>

All outputs in CLoK are protected by Intellectual Property Rights law, including Copyright law. Copyright, IPR and Moral Rights for the works on this site are retained by the individual authors and/or other copyright owners. Terms and conditions for use of this material are defined in the <u>http://clok.uclan.ac.uk/policies/</u>

#### RESEARCH ARTICLE



WILEY

# Bayesian estimation of heterogeneous environments from animal movement data

# Svetlana V. Tishkovskaya<sup>1</sup> | Paul G. Blackwell<sup>2</sup>

Accepted: 24 March 2021

<sup>1</sup>Faculty of Health and Wellbeing, University of Central Lancashire, Preston, UK

<sup>2</sup>School of Mathematics and Statistics, University of Sheffield, Sheffield, UK

#### Correspondence

Svetlana V. Tishkovskaya, University of Central Lancashire, Preston PR1 2HE, UK. Email: s.tishkovskaya@gmail.com

#### **Funding information**

National Centre for Statistical Ecology, Grant/Award Number: EPSRC/NERC grant EP/1000917/1; Natural Environment Research Council

#### Abstract

We describe a flexible class of stochastic models that aim to capture key features of realistic patterns of animal movements observed in radio-tracking and global positioning system telemetry studies. In the model, movements are represented as a diffusion-based process evolving differently in heterogeneous regions. In this article, we extend the process of inference for heterogeneous movement models to the case in which boundaries of habitat regions are unknown and need to be estimated. Data augmentation is used in reconstructing the partition of the heterogeneous environment. The augmentation helps to diminish the impact of uncertainty about when and where the animal crosses habitat boundaries, and allows the extraction of additional information from the given observations. The approach to inference is Bayesian, using Markov chain Monte Carlo methods, allowing us to estimate both the parameters of the diffusion processes and the unknown boundaries. The suggested methodology is illustrated on simulated data and applied to real movement data from a radio-tracking experiment on ibex. Some model checking and model choice issues are also discussed.

#### **KEYWORDS**

animal movement, data augmentation, diffusion, heterogeneous environment, Markov chain Monte Carlo, Ornstein–Uhlenbeck process

## **1** | INTRODUCTION

Recent technological advances in animal tracking systems have made complex, spatially explicit, high-frequency datasets of wildlife behavior and movement available, such as those derived from the global positioning system (GPS). GPS telemetry technology provides nearly continuous, systematically scheduled datasets of locations that allow the details of an animal's movement to be characterized and related to its environment.

This article contributes to the emerging data-rich field of statistical methods for the analysis of individual movements. The availability of new data sources provides opportunities for a more complete and complex analysis of movement processes than was previously possible. Developments in tracking technologies have advanced the study of animal movement and motivated the development of new theoretical frameworks and novel data analysis tools (Hooten et al., 2017). In particular, models formulated in continuous-time (see, e.g., Blackwell, 1997; Brillinger et al., 2001; Dunn & Gipson, 1977) are important because of their flexibility in dealing with irregular or missing data, which can be exploited to develop more efficient designs for telemetry studies (Eisenhauer & Hanks, 2020). Different approaches include building models. This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

 ${\ensuremath{\mathbb C}}$  2021 The Authors. Environmetrics published by John Wiley & Sons Ltd. <sup>2 of 17 |</sup> WILEY

based on tractable processes (Blackwell, 1997; Dunn & Gipson, 1977; see Section 2.1 for details) or using more flexible stochastic differential equation models based on potential functions (Brillinger et al., 2001; Preisler et al., 2004). Models and methods developed for the analysis of animal movement may also have other environmental applications, as in for example, the monitoring of fishing vessels to assess their movement and activity (Gloaguen et al., 2015).

The goal of this article is to develop a continuous-time framework for analysis of individual animal movements combined with estimating the structure of the spatial heterogeneity underlying those movements, using GPS-based or similar location data.

Animal ecology is fundamentally spatiotemporal, with ecological processes occurring on heterogeneous landscapes (Cagnacci et al., 2010). Spatial heterogeneity may occur because of external physical differences across a landscape, for example, different habitat types, or because of internal behavioral responses to location that reflect memory, social interactions, or some kind of strategy in the use of space, for example, a long-term home range or territory. Environmental heterogeneity is built into our movement model by assuming that the landscape consists of a finite set of regions and that the parameters of the movement model are specific to each region. A region may represent an actual "patch" of a habitat type, the presence or extent of which is not known directly, or an area that the animal regards as meaningful but which cannot be observed directly, such as its territory or the core area of its home range. We propose a modeling approach in which location data over time are used to estimate not only the movement parameters but also the unobserved bound-aries between regions. Inference for all parameters of interest is performed within the Bayesian approach, because the model used has a complicated structure and because some prior information from ecologists is often available.

Our approach estimates unobserved boundaries from movement data, and so is much more general than the typical case where regions are known, for example, from independent mapping of habitat features, and only the movement parameters associated with them need be estimated (see Blackwell et al., 2016; Harris & Blackwell, 2013). We do however assume that this structured heterogeneity is constant over time; in contrast, Scharf et al. (2019) look at the case where the heterogeneity takes the form of attraction to a physical feature that varies over time but can be directly observed, while Wang et al. (2019) consider the interaction of spatial patchiness and time-varying resources due to depletion by the observed animals. Our approach imposes a more specific spatial structure than is found in, for example, the potential-based models mentioned above; however, our aim is that this structure has biological meaning, increasing the interpretability of the models.

In the case where the heterogeneity is an internal behavioral feature, our approach shares its aims with methods for identifying the core area of an animal's home range that are based on independent observations of location. For example, Wilson et al. (2010) describe a method for constructing the "core area" based on observations forming a Poisson process with a higher intensity inside the core than outside. Their boundary construction is based on kernel density estimation, rather than the parametric form that we use, but they limit their scope to locations that are independent over time, rather than modeling autocorrelated movement.

Our article is organized as follows. Section 2 describes a diffusion model of movement incorporating spatial heterogeneity. Given unknown boundaries, it is natural to consider the observed data as incomplete, lacking observations on the times and locations of transitions across boundaries. Section 3 presents a method for reconstructing the animal paths between observed locations using data augmentation. In Section 4 we describe in more detail the implementation of spatial heterogeneity used in the article. Markov chain Monte Carlo methods are used to fit the model to data and Section 5 briefly examines some aspects of our implementation. We illustrate the methodology with two examples, with simulated data in Section 6 and with GPS data in Section 7, using telemetry data on animal movements to classify the landscape into two regions and estimate the movement parameters within each. The article concludes with a discussion of possible extensions of our method in Section 8.

## 2 | MOVEMENT MODEL

### 2.1 | Modeling movements in a heterogeneous environment

A useful and popular conceptual model for animal movements is a diffusion, a process in continuous time with continuous sample paths that is the solution to a stochastic differential equation. Let  $X(s) = (X_1(s), X_2(s))$  be a diffusion describing the location of animal at the time *s* (see, e.g., Patterson et al., 2016). We consider a flexible class of stochastic models constructed from two-dimensional Ornstein–Uhlenbeck (OU) processes. Recall that X(s) follows an OU process with parameters ( $\mu$ , B,  $\Lambda$ ) if

$$X(s+\tau)|\{X(s)=x_s\} \sim \mathcal{N}\left(\mu + e^{B\tau}(x_s-\mu), \Lambda - e^{B\tau}\Lambda e^{B'\tau}\right).$$
(1)

WILEY 3 of 17

The OU process serves as a model for animal movement paths in its own right, taking into account high dependency between successive observations. It was first proposed by Dunn and Gipson (1977) and is now a well-established tool for conducting statistical analyses of animal movement data at the individual level (Hooten et al., 2017; Patterson et al., 2016). The OU process has proved to be a workable model for describing movement with physically interpretable parameters; see, for example, Worton (1995) and the review paper of Johnson et al. (2008). This model describes a movement pattern having a stationary distribution with elliptical contours, with the locations concentrated in the center. The parameter  $\mu$  represents this central attraction point which can have various biological interpretations, such as food or conspecific attraction. The parameter  $\Lambda$  is a positive-definite matrix that determines the spread of locations at which the animal is found. The matrix *B* must be stable; Blackwell (2003) notes that the case in which  $B = \beta I$  is a diagonal matrix for some scalar  $\beta < 0$  is usually sufficient in practice.

An important feature of a diffusion model is that it can incorporate observations that are irregular in time. Though many telemetry observations are made at equally spaced time points, the ability of the model to handle arbitrary observation times is useful in many situations, including accounting for missing observations without need for imputation, allowing comparison or synthesis of experiments on different regular timescales, or dealing with studies where sampling intervals vary by design. See Harris and Blackwell (2013), Blackwell et al. (2016), and Eisenhauer and Hanks (2020) for examples and further discussion.

While the limitations of this model are well recognized, in many cases it is a useful approximation under which inferences can be made about the parameters of the home range. Dunn and Gipson's model was extended by Blackwell (1997). In Blackwell's model, for which inference was described in detail in Blackwell (2003), an individual switches between different OU processes, representing different behaviors. A process in which this switching occurs is called a mixed diffusion.

Animals move differently in different regions, and complex patterns of behavior are driven not only by biological mechanisms but also by spatial heterogeneity. Harris (2007) and Harris and Blackwell (2013) considered an extended mixed diffusion model that incorporates spatial heterogeneity, with the processes that describe movement changing according to location as well as behavioral state. We consider this heterogeneous case; we restrict our attention to the single-behavior model, but the extension to the multibehavioral setting is straightforward. We will focus on the study of the spatial heterogeneity, in particular, how it can influence the movement process, and how in turn the movement can help us learn about the heterogeneity, whether internal or external in the sense described in Section 1. Note that this takes us well beyond the class of "separable" models (Harris & Blackwell, 2013) with known spatial structure, for which inference is addressed by Blackwell et al. (2016). We consider the heterogeneous environment  $\mathcal{R}$  as a partition containing a finite set of nonintersecting regions  $R_i$ , with  $\mathcal{R} = \bigcup_{i=1}^k R_i$ . A single pattern of movement is associated with each region, so that the overall diffusion evolves as an OU process with parameters ( $\mu_i, B_i, \Lambda_i$ ) when in region  $R_i, i = 1, \ldots, k$ .

Since we are considering only a finite set of regions, and therefore of movement patterns, any model of this kind is directly related to a model that switches randomly (as in Blackwell, 2003) between those same movement patterns, differing only in that the "switches" are explained by the spatial heterogeneity rather than forming a completely random realization of some underlying Markov chain of "behaviors" that is homogeneous in both space and time. The spatially heterogeneous model, where applicable, will therefore have more explanatory and predictive power than the homogeneous one.

The goal of this article is to develop a method for inference from animal movement that allows use of the model described above in the case where the boundaries between the regions are unknown and have to be estimated from the data on movements. In terms of implementation, this extension adds some extra unknown parameters to those that govern movement. If  $\Psi_i$  is a suitable parametric representation of the boundaries for region  $R_i$ , then altogether the following set of parameters is to be estimated:  $\Theta = (\Theta_1, \ldots, \Theta_k)$ , with  $\Theta_i = (\Omega_i; \Psi_i)$  where  $\Omega_i = (\mu_i, B_i, \Lambda_i)$  are diffusion parameters for region  $R_i$ .

### 2.2 | Likelihood

Let  $x_t = (x_{1t}, x_{2t})$  be the location in two-dimensional space  $\mathbb{R}^2$  of an animal at time index *t*. Then all observations of the animal's locations are represented by  $\mathbf{X} = (x_0, x_1, \dots, x_n)$ . For most of the paper we assume that observations are made at

4 of 17 WILEY

equally spaced time points so that  $\tau$  in expression (1) is constant and observation times are an integer multiple of the constant sampling interval. However, it is important that the model can straightforwardly handle arbitrary observation times.

The inference for parameter  $\Theta$  is performed within a Bayesian framework, using the Markov chain Monte Carlo (MCMC) method. Implementation of the MCMC algorithm is described in Section 5. Here we discuss some assumptions and approximations involved in calculation of the likelihood

$$L(\Theta|X) = p(x_0|\Theta) \prod_{t=1}^{n} p(x_t|x_{t-1},\Theta).$$
<sup>(2)</sup>

Firstly, the likelihood depends on the assumptions made about the initial observation  $x_0$ . There are several practical approaches to evaluating the marginal density function for the first location  $p(x_0|\Theta)$ . Johnson et al. (2008) suggest that, if the movement model  $p(x_t|x_{t-1}, \Theta)$  is stationary, then  $p(x_0|\Theta)$  can be approximated by  $p^{\infty}(x_0, \Theta)$ , the temporal limiting distribution of the movement model, or its weighted version  $w(x_0)p^{\infty}(x_0, \Theta)$  where the weighting function w(x) depends on the selection of resources available for the animal. Since the OU process is stationary, one can use its limiting (equilibrium) distribution  $N(\mu, \Lambda)$  for calculating  $p(x_0|\Theta)$ . In the mixed case, Blackwell (2003) proposes to extract the information contained in  $x_0$  by generating the part of the trajectory before time t = 0. In this reverse simulation, the path  $x_{-s}, \ldots, x_{-1}$  is reconstructed and the extra term  $p(x_0|x_{-1}, \Theta)$  can be used in the likelihood. Another approach is to condition on  $x_0$  and simply remove the  $p(x_0|\Theta)$  term in (2) (Christ, ver Hoef, & Zimmerman, 2008; Johnson et al., 2008). If the time series is long, the loss of information in this approach is small but the complexity of the likelihood is greatly reduced. We use the latter approach because the number of observations in modern studies is usually large.

The conditional distribution of the animal's position  $x_t$  given its position  $x_{t-1}$ , when both locations are in the same region  $R_l$ , is given by the OU process with parameters  $\Omega_l = (\mu_l, B_l, \Lambda_l)$ , making the assumption that the process does not exit the region and return to it between observations, in other words, that there are no unobserved sojourns to different regions. While in some cases this may seem like an oversimplification, it is likely to be accurate when observations are close together in time, either because the actual data are collected at high frequency or given the data augmentation described below.

When  $x_{t-1}$  and  $x_t$  lie in different regions  $R_{l_1}$  and  $R_{l_2}$ , we need to make some assumptions to calculate the transition density for the boundary-crossing move from  $x_{t-1}$  to  $x_t$ . Since the movement processes at the start and end of the interval are known, it seems natural to approximate the conditional distribution  $p(x_t|x_{t-1}, \Theta)$  by a mixture of these two processes. That is, if  $p_{OU}(x_t|x_{t-1}, \Omega_{l_1})$  denotes the conditional density of the OU process in the region  $R_{l_1}$ , then the transition density on an interval including a boundary crossing is approximated by

$$p(x_t|x_{t-1},\Theta) = w_{l,l_2} \times p_{OU}(x_t|x_{t-1},\Omega_{l_1}) + (1 - w_{l,l_2}) \times p_{OU}(x_t|x_{t-1},\Omega_{l_2}),$$
(3)

where  $l_1, l_2$  are defined by  $x_{t-1} \in R_{l_1}$  and  $x_t \in R_{l_2}$ . The adequacy of this approximation depends crucially on the choice of the coefficient  $w_{l_1 l_2}$  for weighted distributions in (3) and on the length of the time interval between observations. The movement process near the boundary may be complicated; the different habitat types in the vicinity of an animal may influence where and how that animal will move. This means that the optimal choice of weight  $w_{l_1 l_2}$  would depend both on the regions adjoining the boundary and also on the animal's precise location. Choosing this coefficient is not straightforward in the case of unknown boundaries between regions. On the other hand, the contributions of both of the component OU processes to the likelihood can be very important when number of transitions across the boundaries is large.

Here we ensure that the limitations of this approximation do not have a large effect, using an approach based on reconstructing intermediate points on the animal's path. When the animal crosses a boundary, the term  $p(x_t|x_{t-1}, \Theta)$  will contribute to the overall likelihood some amount of uncertainty due to the approximation (3), and our reconstruction process aims to reduce this effect by ensuring that only a short time interval is affected. As a result, the approximation for the transition density is less influenced by the coefficients for the weighted distributions in (3), and the weight  $w_{l_1l_2}$  can be chosen flexibly. This approach also reduces the effect of our assumption of 'no unobserved sojourns' which means that all visits to regions are represented by locations; we need only assume that this holds for the augmented paths, not for the actual data. The details of the animal path reconstruction are given in the next section.

## **3** | RECONSTRUCTING TRANSITIONS BETWEEN REGIONS

Information about the times and locations of transitions across boundaries is generally not available. So when the animal's locations at successive time points fall into different regions, calculation of the conditional distribution  $p(x_t|x_{t-1}, \Theta)$  necessarily involves some approximation. Our approach is to reconstruct the animal paths at instants between the observed points. Assume x(s) and  $x(s + \tau)$  are two successive animal locations measured at times s and  $s + \tau$ . Using the proposed movement model and information contained in the observed locations, we generate locations at intermediate time points  $s < s_1, \ldots, s_p < s + \tau$ . These time intervals may be equally spaced or they may become more fine in the vicinity of a boundary. These simulated points  $z(s_1), \ldots, z(s_p)$  are used as a reconstruction of the animal path between the two measurements. If the animal crosses a boundary between points  $z(s_k)$  and  $z(s_{k+1})$ , say, then in the reconstructed trajectory

 $x(s), z(s_1), \ldots, z(s_k), z(s_k + \delta), \ldots, z(s_p), x(s + \tau)$ 

the approximation of the transition on interval ( $s, s + \tau$ ) is replaced by the approximation of transition on the much shorter interval ( $s_k, s_k + \delta$ ). As a result, it helps to diminish the impact of uncertainty about boundary crossings on the likelihood and potentially to extract some additional information from the observations.

Unknown animal locations and times of boundary transitions can be considered within the context of measurement imprecision. Frair et al. (2010) in their review note that there are different approaches to reduce the effects of location imprecision including differential corrections, interpolating and smoothing animal paths. Our approach is to interpolate the collected locations by reconstructing animal movements between fixes, sampling repeatedly to allow propagation of uncertainty.

## 3.1 | Data augmentation

In reconstructing animal paths between the regions of a heterogeneous environment, we use the data augmentation technique. The observed data X are treated as incomplete and are augmented by latent (unobserved) data Z using a method similar to that suggested by Tanner and Wong (1987). This approach is common in movement modeling, either to add the times (and perhaps locations) of behavioral switches in between observations, for exact fitting of continuous-time models (e.g., see Blackwell, 2003; Blackwell et al., 2016), or to add "true" locations at the times of the existing observations when allowing for measurement error (e.g., Johnson et al., 2008). In those examples, and in the present work, the augmentation of the data takes place in conjunction with inference about the model parameters, with each informing the other. Thus the sampled parameters directly incorporate all the information from the data. This contrasts with a two-stage approach where refinement of the trajectories is carried out first, using a simpler movement model, and then inference from the model of interest is carried out separately based on those refined paths. A two-stage approach may be appropriate in some circumstances; see McClintock (2017) and Scharf et al. (2017). for details and discussion. In the present case, a two-stage approach is unlikely to perform well because uncertainty in the boundaries has a large effect on reconstructed paths. Our approach of simultaneous augmentation and parameter inference allows us to extract additional information about the movement of the animal between regions; this could be combined with augmentation to address behavior or observation error, but we omit those factors here for simplicity. Mathematically, the augmentation method makes it possible to express the required posterior density in the form

$$p(\Theta|x) = \int_{Z} p(\Theta|z, x) p(z|x) dz.$$
(4)

WILEY 5 of 17

To evaluate this integral, we draw posterior simulations of the joint vector of unknowns ( $\mathbf{Z}, \Theta$ ) and then focus on the estimands of interest. The process is implemented in the following iterative scheme where the augmented data  $\mathbf{Z}$  are treated as variables.

**Step 1.** Simulate  $\Theta \sim p(\Theta|z, x)$ . Accept or reject the simulated  $\Theta$ . **Step 2.** Using  $\Theta$  from Step 1 simulate  $\mathbf{Z} \sim p(z|x, \Theta)$ . Accept or reject the simulated  $\mathbf{Z}$ .

Accept–reject steps are done by a single-component Metropolis–Hastings algorithm for  $\Theta$  and an independence sampler for **Z**. Realization of these steps will be described in Section 5. Iterative sampling results in the set of observed data **X** being expanded by adding augmented data **Z** and the full data set is now  $\mathbf{Y}^{aug} = \mathbf{X} \cup \mathbf{Z}$ .

For the imputation Step 2 above we need a method for proposing the latent variable **Z**. We propose a path between two successive locations  $x_{t-1}$  and  $x_t$  from a Wiener process (Brownian motion) tied to the observations at the endpoints, that is a Brownian bridge. This process appears to be a sufficiently good representation of the animal movement paths

for this purpose, over short intervals. In practice we use a time-discretization approach, proposing values of the path at specific intermediate times using this Brownian bridge.

Samples from the Brownian bridge could be generated sequentially or simultaneously, but, for convenience, the augmentation is carried out recursively. For each pair of observations  $x_t$  and  $x_{t+1}$  we form a Brownian bridge  $\mathcal{B}^t(s)$ , t < s < t + 1 with  $\mathcal{B}^t(t) = x_t$ ,  $\mathcal{B}^t(t+1) = x_{t+1}$ , and generate its value at the mid-point of the time interval,

$$\mathcal{B}^t\left(t+\frac{1}{2}\right) \sim \mathrm{N}\left(\frac{1}{2}(\mathcal{B}^t(t)+\mathcal{B}^t(t+1)),\frac{1}{4}\sigma^2\right),$$

where the variance  $\sigma^2$  is derived from the variances of OU processes in the corresponding regions. This bisection of the time intervals is then repeated by constructing a Brownian bridge  $\{B_j^t(s)\}, s \in (tj', tj'')$  on each of the intervals (tj', tj'') formed, and generating the value for each one at *its* mid-point,

$$\mathcal{B}_j^t\left(\frac{t_j'+t_j''}{2}\right) \sim \mathrm{N}\left(\frac{1}{2}(\mathcal{B}^t(t_j')+\mathcal{B}^t(t_j'')),\frac{1}{4}(t_j''-t_j')\sigma^2\right).$$

This refining algorithm fills in the  $p_t$  required intermediate points between  $x_t$  and  $x_{t+1}$ ; we generally take each  $p_t$  to be of the form  $2^i - 1$  (i = 1, 2, 3, ...). The (time-ordered) points generated in this way form the augmented data  $z_{tj}$ ,  $j = 1, ..., p_t$ .

For computational stability, unobserved trajectories are reconstructed between all animal locations in **X**. For each observed pair  $x_t$  and  $x_{t+1}$  the values of latent data  $z_{t,1}, z_{t,2}, \ldots, z_{t,p_t}$  are generated. We then have the extended data set arrayed in the form

| $x_0$ ,    | $Z_{01},$    |     | ••• | $Z_{0p_0}$        |
|------------|--------------|-----|-----|-------------------|
| $x_1$ ,    | $z_{11},$    |     | ••• | $Z_{1p_1}$        |
| •••        | •••          | ••• |     | •••               |
| $x_{n-1},$ | $Z_{n-1,1},$ |     |     | $Z_{n-1,p_{n-1}}$ |
| $x_n$ .    |              |     |     |                   |

#### 3.2 | Augmented likelihood

The above scheme provides us with augmented data  $\mathbf{Y}^{aug} = (y_1, \dots, y_T)$  which includes both observed **X** and unobserved **Z** components as described in Section 3.1. The "complete-data" augmented likelihood in this notation is

$$L^{aug}(\Theta|\mathbf{Y}^{aug}) \propto \prod_{t=1}^{T} p(y_t|y_{t-1},\Theta),$$
(5)

where the "complete-data" conditional movement density is given by

$$p(y_t|y_{t-1},\Theta) = \begin{cases} p_{OU}(y_t|y_{t-1},\Omega_l), & y_{t-1}, y_t \in R_l(\Psi_l) \\ w \cdot p_{OU}(y_t|y_{t-1},\Omega_{l_1}) \\ + (1-w) \cdot p_{OU}(y_t|y_{t-1},\Omega_{l_2}), & y_{t-1} \in R_{l_1}(\Psi_{l_1}), y_t \in R_{l_2}(\Psi_{l_2}), l_1 \neq l_2. \end{cases}$$
(6)

In the above mixture weights *w* are taken to be equal across all regions. The finer time-discretization, resulting from the augmentation, makes this approximation acceptable. In practice, we use w = 0.5, assuming that the two OU processes make equal contributions to the movement density in the vicinity of a boundary crossing.

### **4** | IMPLEMENTING THE MODEL

While the described method is applicable to any finite partition of space, for illustrative purposes we shall restrict ourselves to a simple example. We consider a simplistic study area that comprises two regions: an inner region with a circular

6 of 17

boundary where the animal spends most of its time and an outer zone of relatively less intense use. Such a circular model can be seen as an approximation for a patch of habitat that forms a convex region and is treated as homogeneous due to the uniformity of the movement pattern when compared with the whole study area. More generally, it may be a simplified representation of the animal's internal idea of its territory or the core of its home range. This illustrative case is supported by considerations of efficiency due to shape, in a number of contexts. Summarizing properties of different types of landscape patches, Barnes (2000) points out that a circular patch minimizes the amount of edge compared, for example, to a thin, rectangular strip, which has only a narrow band of interior habitat. One example of circular habitat usage is provided by red squirrels. These animals are considered to be territorial in coniferous forests. The article by Gurnell (1984) examining home range and territoriality in red squirrels has indicated that all the observed squirrels patrolled and defended areas of circular or oval shape, centered on large caches and a nest site. Similarly, biological studies suggest (see, e.g., Linn et al., 2007) that the home ranges of elephant-shrews are compact and quite symmetrical, with an inner zone of intense use surrounded by a zone of relatively less intense use and a variable outer zone of sparse use. The symmetrical perimeter of the intensively used area can be approximated by a circle. Space use within this perimeter is strongly symmetrical, with the most intense use at the center. These analyses of range use provide motivation for use of the two-region OU model with a circular boundary to model both movement and use of space. A simplistic circular model is also used as an approximation of the home range for some large carnivores, for example mountain lions (cougars), that have large territories, usually oval or circular in shape (Russell et al., 2012). Finally, if the movement process in the outer region is taken to be Brownian motion (available as a limiting case of the OU process), the radius of the circle can represent the (unknown) maximum range of attraction of the central food source, and so on; such a model will be transient rather than stationary.

Modeling examples of an idealized animal's territory or home range represented by a circle are given in Harris and Blackwell (2013). They provide two simulated examples of a circular region: uniform environment with a central point of attraction within the circle and unused patch of habitat within a home range with a central point of repulsion inside the circle.

# 5 | IMPLEMENTATION OF THE MCMC ALGORITHM

Often the time interval  $\tau$  between collected radio-tracking data is a constant, and thus it is more convenient to use an alternative parameterization of the OU process where  $\tau$  is fixed. We use the same reparameterization for diffusion parameters as in Blackwell (2003) and Harris (2007) with  $\Lambda$  and *B* matrices replaced by

$$\Gamma = \exp(B\tau), \quad \Phi = \Lambda - \Gamma\Lambda\Gamma',$$

where  $\exp(\cdot)$  is the matrix exponential and parameter  $\Phi$  is the covariance matrix of an observation from an OU process conditional on an observation at a time  $\tau$  earlier. To incorporate irregular data,  $\tau$  can be taken to correspond to the most frequently occurring interval between observations. This parameterization is more convenient for MCMC simulation and gives us the following vector of unknown parameters for region  $R_i$ :

$$\Theta_i = (\mu_i, \Gamma_i, \Phi_i; \Psi_i), \quad i = 1, \dots, k.$$

In the remainder of the article we use both parameterizations depending on the context.

As outlined in Section 3.1, after appropriate initialization of **Z** and  $\Theta$  the MCMC algorithm alternates between the following steps.

#### Step 1. Sampling $\Theta|Z, X$ (random-walk Metropolis–Hastings)

For each parameter  $\theta \in \Theta$ :

propose  $\theta'$  from a proposal distribution centered on  $\theta$ ;

accept or reject  $\theta'$  using a Hastings ratio based on the augmented likelihood (5).

## Step 2. Sampling $Z|\Theta, X$ (independence sampler)

For *t* in 0, ..., n - 1:

propose  $z'_{t,1}, \ldots, z'_{t,p_t} \sim B^t | x_t, x_{t+1}$  from the Brownian bridge of Section 3.1, independently of  $z_{t,1}, \ldots, z_{t,p_t}$ ; accept or reject  $z'_{t,1}, \ldots, z'_{t,p_t}$  using a Hastings ratio based on appropriate terms of the augmented likelihood (5).

For our examples, we define the parameterization of boundaries  $\Psi_i$  by fixing a circular boundary between two regions so that the study area consists of a circular "core"  $R_1$  and a region  $R_2 = \mathcal{R} \setminus R_1$ , outside that core. Since our examples are intended to model cases with the most intense use at the center, we assume that two regions have the same attraction point and both OU processes share a common center:  $\mu_1 = \mu_2 = \mu$ .

Blackwell (2003) notes that in practice the isotropic case is usually sufficient, with the matrix *B* given by  $\beta I$  where  $\beta < 0$  is some scalar and *I* is the identity matrix. The isotropic assumption gives us  $\Gamma = \gamma I$  with  $0 < \gamma < 1$ . Also, since we model the circular home range, it is biologically natural to assume that matrix  $\Lambda = \lambda I$  and an attraction point inside the circle implies  $\lambda_1 > \lambda_2$ . Parameter  $\Phi$ , in turn, is reduced to a scalar  $\phi = \lambda(1 - \gamma^2) > 0$  with  $\Phi = \phi I$ . This parameterization gives us the vector of unknown parameters  $\Theta = (\mu, \gamma_1, \gamma_2, \phi_1, \phi_2, \rho)$  where  $\rho$  is radius of the circular boundary and  $\mu \in \mathbb{R}^2$  is the position of the attraction point.

In the examples, we use vague proper priors for all parameters, assuming that there is little prior information about parameters. Details of the priors are given in Web Appendix A, with prior and proposal distributions for each parameter used in the MCMC simulation summarized in Web Table 1. The parameters of the proposal distributions are chosen to give Metropolis–Hastings steps with an acceptance rate that is not too far from optimal values.

To illustrate the described approach and its performance, in the next sections we will consider two examples, one with simulated data (Section 6) and the other with radio-tracking data on movements of an ibex (Section 7).

#### **6** | SIMULATION EXPERIMENT

To demonstrate the benefits of using the augmentation for estimation of parameters, we provide the results of a simulation experiment. In the experiment we compared three computational schemes: estimation of the parameter  $\Theta$  from the original observed points and from augmented data sets with either one or three points imputed between each pair of observations (schemes  $A_0$ ,  $A_1$ , and  $A_3$ , respectively).

To illustrate the possible gains from using augmentation, we conducted the simulation experiment for a small sample size for which parameter estimation seems more challenging. Animal paths of n = 100 observations were simulated from mixed OU processes, representing movement patterns in a study area consisting of two regions, a circular core and an outer zone. The animal has an attraction toward a central core area but also makes excursions from this area. The same parameters were used throughout all simulations, and we set  $\Theta_T = (\rho = 2, \mu = (0, 0), \beta_1 = -3, \beta_2 = -0.5, \lambda_1 = 1, \lambda_2 = 0.8)$ . Structurally, the model is the same as that fitted to the ibex data; the parameter values are broadly comparable, but are set to represent a somewhat different case, where the function of the movement parameters that is most easily estimated directly from the data,  $\phi_i$ , varies between regions much less than is the case for the ibex, making the estimation more challenging. (In the parameterization used in the MCMC algorithm, the simulations have  $\gamma_1 = 0.050$ ,  $\gamma_2 = 0.61$ ,  $\phi_1 = 1.00$ ,  $\phi_2 = 0.506$ .) While these parameter choices are not meant to mimic a particular species, they plausibly represent observations on a similar time scale to the ibex data, around every 4 h.

The parameter estimates for each path were calculated using MCMC sampling. The number of MCMC iterations in simulations  $A_0$ ,  $A_1$  was taken as N = 10,000; for scheme  $A_3$  we used a slightly greater number of iterations, N = 12,000. These chains, after the first burn-in steps were discarded, were used for calculating posterior estimates  $\hat{\Theta}$ . Chains obtained with augmented datasets tend to converge more slowly and the number of burn-in iterations was adjusted for each scheme to take into account convergence rate.

For each simulated path we therefore had three sets of parameter estimates, for  $A_0$ ,  $A_1$ , and  $A_3$ . Parameter estimates were compared using a root mean square error (RMSE) statistic with the distance  $d_{\Theta} = d(\hat{\Theta}, \Theta_T)$  between the parameter estimate  $\hat{\Theta}$  and the true value of parameter  $\Theta_T$  measured as Euclidean distance:

$$\text{RMSE}_{\Theta} = \left[ E(d(\widehat{\Theta}, \Theta_T))^2 \right]^{1/2}$$

In the simulation experiment we were interested to examine to what degree the algorithm with data augmentation allows us to make better inference concerning unknown boundaries and movement parameters. To explore this, we considered separately the distances  $d_{\Psi}$  and  $d_{\Omega}$  for boundary parameters  $\Psi$  and diffusion parameters  $\Omega$  accordingly. Distances at the level of individual parameters are not given, because of the correlations between parameters within each of  $\Psi$  and  $\Omega$ .

The results of the parameter estimate comparisons for 50 simulated animal paths are presented in Table 1. The table shows the values of statistic RMSE<sub>*P*</sub><sup>*A<sub>i</sub>*</sup> for each computational scheme *A<sub>i</sub>* and each parameter *P* =  $\Theta$ ,  $\Psi$  and  $\Omega$ , pooled across

**TABLE 1** RMSE for parameter  $\Theta$ , boundary parameters  $\Psi$ , and diffusion parameters  $\Omega$  with percentage gain in decreasing RMSE due to augmentation in brackets

|                        | $A_0$ | $A_1$         | $A_3$         |
|------------------------|-------|---------------|---------------|
| $\text{RMSE}_{\Theta}$ | 1.12  | 0.74 (34.04%) | 0.65 (41.75%) |
| $RMSE_{\Psi}$          | 0.84  | 0.39 (52.70%) | 0.51 (38.95%) |
| $\text{RMSE}_\Omega$   | 0.72  | 0.61 (15.13%) | 0.44 (38.84%) |

samples. In brackets we show the relative reduction in RMSE due to augmentation,  $\left(\text{RMSE}_{p}^{A_{0}} - \text{RMSE}_{p}^{A_{i}}\right)/\text{RMSE}_{p}^{A_{0}}$ , i = 1, 3, expressed as a percentage. Because the sample distribution of RMSE statistics was skewed for schemes with augmentation, the average RMSE values are presented in the table as estimates of medians. Distribution of the statistics RMSE\_{\Theta}^{A\_{i}}, RMSE\_{ $\Psi}^{A_{i}}$ , and RMSE\_{\Omega}^{A\_{i}} summarized using box plots are shown on Web Figures 1, 2, and 3, respectively, for each computational algorithm. It should be noted that the outliers in Web Figures 1 and 3 for the  $A_{3}$  scheme reflect the fact that runs with more augmentation required more iterations to achieve convergence for some of randomly generated 'animal paths', in particular for diffusion parameters. This would have to be taken into account when implementing the method on practice.

The simulation studies show definite changes in estimating both the boundary and OU movement parameters. In terms of RMSE, even a small amount of augmentation gave clear improvements in the parameter estimation. Table 1 shows that for RMSE<sub> $\Theta$ </sub>, there is a 34% decrease with one augmented point and a 42% decrease with three augmented points between each observation, compared with no augmentation scheme. The experiment demonstrated that augmentation is beneficial when an animal path includes a large number of boundary crossings. The effect of using augmentation is apparent even with a small number of added points between each animal location, and is particularly pronounced for boundary parameters  $\Psi$ ; RMSE statistic RMSE<sup>A<sub>1</sub></sup> represents a 53% improvement in boundary estimates compared with the  $A_0$  scheme, though the improvement does not progress with further increasing number of augmented points which is, probably, related to the difference in the number of potential boundary crossing between  $A_1$  and  $A_3$ .

In the web-based Supplementary Materials for the article, Web Figure 4 provides a visual illustration of the estimation of  $\Theta$  for the  $A_0$  and  $A_3$  schemes for one arbitrarily chosen simulated path. The example demonstrates that the augmentation gave a visible improvement in the selection of the boundary between regions.

### 7 | IBEX EXAMPLE

#### 7.1 | Data and model

The Ibex dataset contains the GPS relocations of an ibex during 15 days in the Belledonne mountain range (French Alps). Data on relocations were collected roughly every 4 h. During the study period 71 observations were recorded. The data can be found in the data repository of the package adehabitat Calenge et al. (2009); the original source of the data is "Office national de la chasse et de la faune sauvage," France.

Ibex locations are shown in Figure 1, with the initial observation marked by a triangle. The figure suggests that the animal movements have a centralizing tendency toward a core area at the upper part of the plot, which can be interpreted as a foraging home range, but the animal also made some short excursions from this area. The foraging area has a center of attraction which might be, for example, a concentrated food source. Animal movements around the attraction can also be interpreted as "area-restricted." We fit the model to the data with two OU processes sharing an attraction point. One process will describe the animal relocations from starting point toward the core region and outside it ("exploratory" movements) and the second will represent short animal movements within circular foraging habitat ("resident" movements). We interpolated each pair of observations with three augmented points.

### 7.2 | Results

Figure 1 shows the estimated boundary between the two regions together with the observations on ibex movements. Table 2 shows the posterior mean estimates of all parameters, together with their standard deviations (SD) and highest posterior density (HPD) intervals. The posterior distribution of each parameter was estimated from a sample of 30,000



| Parameter  | Mean  | SD    | 95% HPD interval |
|------------|-------|-------|------------------|
| ρ          | 0.946 | 0.024 | (0.904, 1.010)   |
| $\mu_1$    | 0.257 | 0.034 | (0.209, 0.341)   |
| $\mu_2$    | 2.113 | 0.025 | (2.072, 2.180)   |
| $\gamma_1$ | 0.243 | 0.044 | (0.160, 0.332)   |
| $\gamma_2$ | 0.884 | 0.038 | (0.812, 0.959)   |
| $\phi_1$   | 0.085 | 0.010 | (0.068, 0.104)   |
| $\phi_2$   | 1.770 | 0.184 | (1.436, 2.147)   |

**TABLE 2** Posterior parameter estimates with their SD and corresponding HPD intervals for ibex data

MCMC iterations, after a burn-in period of 10,000 iterations and meeting convergence criteria. The convergence was inspected using a visual graphical assessment. The diagnostic was performed for each parameter in the model using multiple parallel chains starting with different values. The diagnostic trace plots and posterior density plots were displayed and visually analyzed to check that chains are mixing well and that a consistent estimate of the posterior distribution is reached.

Adding augmentation results in placing a wider boundary than the results obtained from nonaugmented data, with the boundary estimates for nonaugmented sample being  $\rho = 0.66 (0.06)$ ,  $\mu_1 = 0.42 (0.06)$ ,  $\mu_2 = 1.94 (0.06)$  (presented as mean (SD)). It could be noted that applying the augmentation leads to reduction in the posterior SDs indicating an improvement in the parameter estimation. The wider boundary from augmented data means that more points located near the boundary are included in the foraging area. Visual inspection of sample paths suggests that it makes sense since these "near-boundary" points correspond to rather short animal moves. Thus augmenting transitions across the boundary helps to identify the type of location, in particular, whether it is exploration or resident movement.

We conducted a small additional simulation to explore how augmentation helps to learn about unobserved visits to a region. In the ibex dataset there are 21 observed boundary crossings for the estimated with three augmented points

FIGURE 1 Estimation of habitat boundary for ibex data

boundary (30% of all registered movements). For a sample of simulated augmentation paths we had 57 (20%) crossings on average (median) that is average of 36 unobserved animal visits to a region.

The means of the estimators in Table 2 show that a fitted OU process outside the central patch tends to be more dispersed with a weaker attraction toward the core area than the process governing animal movements within the foraging region, as would be expected. Note that the parameter  $\gamma_2$ , which controls the strength of drift toward the center for exploratory movements, is close to 1, implying that the corresponding diffusion process is close to Brownian motion.

### 7.3 | Model checking and comparison

Although the two-region model seems to be natural for the ibex data and to perform well, obviously there are other options in setting up a model that can fit the data. As part of our model checking, we compared the two-region model with a model consisting of one region; these two competing models will be labeled as  $M_2$  and  $M_1$ , respectively. In the one-region model, the study area is not partitioned into subregions and it is assumed that animal movements come from a single underlying OU process with parameters  $\Theta^{M_1} = (\mu, B, \Lambda)$ . Correspondingly, for the two-region model we have  $\Theta^{M_2} = (\Omega_1, \Omega_2, \Psi)$  with  $\Omega_i = (\mu, B_i, \Lambda_i), i = 1, 2$ , and  $\Psi = \rho$ . To compare the fit of the different models we use two approaches, both based on posterior predictive model comparison.

As noted in Section 2.1, we could also fit, say, a spatially homogeneous two-state switching OU model to these data. Such a model may well do a good job of capturing the difference between the initial observations that are relatively far from the center and those that are more clustered, but is less likely to adequately represent the later data which include shorter excursions beyond the core. More importantly, a nonspatial model will be less able to *explain* these variations in movement pattern. So, while we acknowledge that a two-state model would probably be selected over a one-region/one-state model, we do not pursue that case formally.

## 7.3.1 | Deviance information criterion

The deviance information criterion (DIC), introduced by Spiegelhalter et al. (2002) and widely used for comparing models in a Bayesian framework, is given by  $DIC = \overline{D(\Theta)} + p_D$ , where  $D(\Theta) = -2 \log p(X|\Theta)$  and  $p_D = \overline{D(\Theta)} - D(\tilde{\Theta})$ . As a point estimate for  $\Theta$  we use the posterior mean.

Because DIC is used as a model choice criterion in our example, we calculate its value associated with the observed likelihood  $L(\Theta|\mathbf{X})$ , without using augmentation at this step. For both models, the values of *DIC* and  $p_D$  were calculated from MCMC output of posterior simulations for all model parameters.  $\overline{D(\Theta)}$  was estimated by the sample mean of the simulated values of  $D(\Theta)$  and  $D(\tilde{\Theta})$  was estimated by substituting in the sample mean of the simulated values of  $\Theta$ . The values obtained for the one-region model are  $DIC^{M_1} = 394.22$  and  $p_D^{M_1} = 4.034$  and for two-region model they are  $DIC^{M_2} = 263.67$  and  $p_D^{M_2} = 5.52$  (based on 3000 simulations for each model). Note that the  $p_D$  values are roughly equal to the numbers of model parameters. The large difference between  $DIC^{M_1}$  and  $DIC^{M_2}$  suggests that DIC supports the two region model, which represents a better fit to the data.

### 7.3.2 | Posterior predictive probability

Another natural approach is to compare models on the basis of the posterior predictive distributions of some specifically chosen test quantity  $T(X, \Theta)$ . The approach is related to the numerical posterior predictive check method described in Gelman et al. (2004), but here we use it with a focus on comparing discrepancies between models rather than checking model fit; we compare  $T(X^{obs}, \Theta^{M_i})$  under two different models  $M_1$  and  $M_2$ .

Let  $\eta_1$ ,  $\eta_2$ , and  $\eta_3$  be successive data points, ordered by their times, say  $s_1$ ,  $s_2$ , and  $s_3$ , respectively. With fixed end points, the middle point  $\eta_2$  follows the OU bridge  $Q_s = Q(s; \eta_1, \eta_3)$ , that is an OU process conditioned to go from  $\eta_1$  at time  $s_1$  to  $\eta_3$  at time  $s_3$ . The conditional density corresponding to this process can be written as

$$p^{Q}(\eta_{2}|\eta_{1},\eta_{3},\Theta) = \frac{p(\eta_{2}|\eta_{1},\Theta) \times p(\eta_{3}|\eta_{2},\Theta)}{p(\eta_{3}|\eta_{1},\Theta)}.$$
(7)



30

40

Observation number

20

FIGURE 2 Ibex data: Comparing conditional probability densities under two models (× for model  $M_1$ ,  $\circ$  for model  $M_2$ )

The test statistic  $T(X^{obs}, \Theta^{M_j})$  for model comparison is based on the conditional probability (7). We consider the log ratio of conditional probabilities for two models:

60

70

50

$$T(X^{obs}, \Theta^{M_1}, \Theta^{M_2}) = \sum_{i=1}^{n-2} \log \frac{p_i^Q\left(\Theta^{M_2}\right)}{p_i^Q\left(\Theta^{M_1}\right)},$$

where

0

$$p_i^Q(\Theta^{M_j}) = p^Q(\eta_2 = x(s_i)|\eta_1 = x(s_{i-1}), \eta_3 = x(s_{i+1}), \Theta^{M_j}).$$

Values of the probability density (7) were calculated under both models for each observed ibex location. Logarithms of these values are plotted on Figure 2.

The figure also shows discrepancies between log conditional probability densities for each observation (lower plot). The sum of these discrepancies gives the value of our test statistic,  $T(X^{obs}, \Theta^{M_1}, \Theta^{M_2}) = 41.23$ . For most observations, the conditional probability density under the two-region model is higher than under the model with just one region.

To further illustrate this approach, we compared the predictive distribution of a particular point from the dataset under two models. Figure 3 shows part of the observed ibex movement path which we use for comparing models. Three successive points  $(x_{17}, x_{18} \text{ and } x_{19} \text{ in the dataset})$ , connected by arrows in time order, represent animal movement in a backward direction, which is unusual with an attraction point model. We compare the two models by examining the question: which model explains this atypical movement path better?

We now take as our test statistic

$$T(X^{obs}, \Theta^{M_i}) = p^Q(\eta_2 = x(s_2)|\eta_1 = x(s_1), \eta_3 = x(s_3), \Theta^{M_i}), \quad i = 1, 2$$

For the one-region model, we calculate  $p^Q(\eta_2|\eta_1,\eta_3,\Theta^{M_1})$  by substituting the corresponding conditional densities of the OU process with parameters  $\Theta^{M_1} = (\mu, B, \Lambda)$  into formula (7). For two regions,  $p(\eta_2 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_2 | \eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_2 | \eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_1, \Omega_2), p(\eta_2 | \eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_2 | \eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_1, \Omega_2), p(\eta_2 | \eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_2 | \eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_2 | \eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_2 | \eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_2 | \eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_2 | \eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_2 | \eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_2 | \eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_2 | \eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_2 | \eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_2 | \eta_1, \Omega_2), p(\eta_3 | \eta_1, \Theta^{M_2}) = p_{OU}(\eta_3 | \eta_1, \Theta^{M_2})$  $p_{OU}(\eta_3|\eta_1,\Omega_2)$ . Since with fixed  $\eta_1$  and  $\eta_3$  the intermediate point can belong to either region  $R_1$  or  $R_2$ , the remaining term is:

$$p(\eta_3|\eta_2, \Theta^{M_2}) = p_{OU}(\eta_3|\eta_2, \Omega_1) \operatorname{Prob}(\eta_2 \in R_1|\Omega_1, \rho) + p_{OU}(\eta_3|\eta_2, \Omega_2) \operatorname{Prob}(\eta_2 \in R_2|\Omega_2, \rho).$$
(8)

**FIGURE 3** Three successive points from ibex dataset chosen for model comparison



The probabilities that  $\eta_2$  belongs to  $R_1$  or  $R_2$  in (8) are calculated by numerical integration over the corresponding region. Test statistic values  $T(X^{obs}, \Theta^{M_i})$  are computed using posterior simulations  $\Theta_{(1)}^{M_i}$ , ...,  $\Theta_{(L)}^{M_i}$  from a MCMC run for each model. For both models, these simulated values of the logarithms of conditional probabilities for observed animal locations  $x(s_1), x(s_2), x(s_3)$  are shown in Figure 4.

The two chains in the upper plot represent values of  $\log(T(X^{obs}, \Theta_{(l)}^{M_1}))$  and  $\log(T(X^{obs}, \Theta_{(l)}^{M_2}))$  for l = 1, ..., L and for L = 2000 runs. Estimated values of conditional probabilities under model  $M_2$  are consistently higher than under model  $M_1$ . Figure 4 also shows the estimated posterior densities of  $T(X^{obs}, \Theta^{M_i})$ . The estimated posterior means for one- and two-region models, respectively are  $\overline{p^Q}(x(s_2)|x(s_1), x(s_3), \Theta^{M_1}) = 1.635 \cdot 10^{-6}$  and  $\overline{p^Q}(x(s_2)|x(s_1), x(s_3), \Theta^{M_2}) = 1.682 \cdot 10^{-4}$ . Thus it follows from the conducted posterior predictive check that the combination of the observed animal locations shown in Figure 3 is more plausible under model  $M_2$ , which suggests that specific patterns in the data are better captured by the two-region model.

#### 8 | DISCUSSION

In this article, we have presented a method for modeling individual animal movement in a heterogeneous environment, while simultaneously learning about that environment. The approach involves applying diffusion models to account for observed movement and using a data augmentation technique to reconstruct animal locations between collected observations. The novelty in the present approach is that boundaries between spatial heterogeneous regions are included into the model as unknown parameters.

The example analyses presented herein were concerned with animal relocations between just two regions, one of which was circular. These are relatively simple models but their analysis gives an indication of how we might expect the algorithm to perform in more complex applications. Even with that simplification in the partitioning of the study area, obtaining parameter estimates was computationally demanding. As the number of observed locations and the number of augmented points per pair of observations increase, this estimation procedure quickly becomes very expensive. In particular, runs with more augmentation require more iterations as well as increased cost per iteration, and the total computational cost increase faster than linearly in the number of augmented points. Estimating parameters for large study areas with sophisticated heterogeneous regions may become problematic, although our method could benefit from a parallel computing framework which would keep the algorithm computationally feasible. However, our results suggest



Chains (on log scale) and posterior densities of conditional probabilities  $p^Q$  at observed point  $x(s_2) = x_{18}$  from ibex dataset FIGURE 4 for two models. Dashed lines on the plot correspond to model  $M_1$ 

that even low levels of augmentation can give appreciable benefits. Even a single augmented point between observations allows for the range of possible times at which a given boundary crossing may take place, while adding more points allows for unobserved visits to a region to be accommodated, capturing the key omissions from the data as naïvely interpreted.

The general framework presented here can be extended in a number of directions. The first practically important generalization would be to consider models with greater behavioral complexity. In addition to estimating movement and boundary parameters for multiple regions, Bayesian inference could be used to account for different behavioral states in each region. Incorporating movement states into our model would be straightforward; the general approach in simpler settings with spatial homogeneity or with fixed boundaries is described in Blackwell (2003) and Blackwell et al. (2016), respectively.

Another extension is directly concerned with modeling spatial heterogeneity. Although the approach developed in this study can be used to model the movements of animals in environments with various spatial structures, describing complex heterogeneous environments which can influence the movement patterns of organisms will require more sophisticated models for partitioning of the study area. One promising means for developing realistic and flexible models for habitat fragmentation or, say, complex patterns of territoriality is to represent the partition parametrically in the form of a random tessellation, as applied in ecological, environmental and other contexts by Blackwell and Macdonald (2000), Blackwell and Møller (2003), and Pope et al. (2019).

More ambitiously, this approach could be combined with the representation by Niu et al. (2016) and Milner et al. (2021) of collective movement as a point following an OU process in a higher-dimensional space. Different movement in different regions then represents the dependence of behavioral interactions on the geometry of the group at any instant, for example, two animals being within a certain distance of each other is equivalent to the point representing their joint locations being within a particular infinite circular cylinder. Our approach could thus help learn about the range and geometry of such collective movement.

We have focused here on models built from the OU process, because of their tractability and flexibility. Some of their limitations, in terms of realism of movement modeling, could be addressed by instead applying our approach to the *integrated* OU process, often known as the continuous-time correlated random walk (CTCRW) in an animal movement context, in which it is the velocity of the animal and not its position that follows an OU process. Johnson et al. (2008) introduced a spatially and behaviorally homogeneous CTCRW, and Michelot and Blackwell (2019) showed how behavioral switching could be incorporated without time-discretization. Russell et al. (2018) included spatial heterogeneity in the form of a spline-based "motility surface" and applied it to high-frequency laboratory data on ant movement, using Euler–Maruyama numerical approximation. Applying our approach, where data are not sufficiently high frequency that changes between fixes can be ignored, would require augmentation with velocities as well as positions, but since our results show that fairly low levels of augmentation have appreciable benefits, it should be readily achievable. For our particular example, we need the animal's location—not just its velocity—to be stationary, and so the CTCRW is not directly applicable without some additional mechanism being included.

In principle, in the Bayesian framework the number of regions *k* can be treated as a variable and included into the model as unknown parameter. Technically, Bayesian inference in this case can be done using the reversible jump MCMC method or similar, since the number of parameters in  $\Theta = (\Theta_1, \dots, \Theta_k)$  can change at each iteration. Pope et al. (2019) shows how this might work in practice, using a random tessellation, as discussed above, with a varying number of regions.

In this study the number of heterogeneous regions k was chosen in Section 7.3 using model comparison tools. It should be noted that the interpretation of DIC in this setting requires caution and further investigation, though the results in our example were sufficiently clear-cut for this not to be an immediate concern. Apart from taking into account the unobserved data  $\mathbf{Z}$ , the model with unknown boundaries can be interpreted as a missing data model problem in the sense that we do not know which region each observation belongs to. We can formalize identifying the region  $R_j$  that observed location  $x_i$  falls in by introducing the indicator variable  $\mathbf{r} = (r_0, \ldots, r_n)$  where  $r_i = (r_{i1}, \ldots, r_{ik})$  so that  $r_{ij} = 1$  if  $x_i \in R_j$  and  $r_{ij} = 0$  if  $x_i \notin R_j$ . Under the model with unknown boundary parameters  $\Psi_i$  the variable  $\mathbf{r}$  should be treated as missing data.

We have shown that is feasible to use autocorrelated location data, such as is increasingly available from GPS tagging, to learn about the spatial form of heterogeneity that influences an animal's movement behavior. We have also indicated a number of ways in which our approach could extend naturally to accommodate more complex models and ecological questions. The method as it stands is computationally expensive, but our results show that even low levels of data augmentation can result in appreciable improvements in the precision of the model fitting.

#### ACKNOWLEDGMENTS

This work was supported by a Daphne Jackson Fellowship to Svetlana V. Tishkovskaya, funded by the Natural Environment Research Council (NERC). The contribution of both authors was also part-funded by EPSRC/NERC grant EP/1000917/1 (National Centre for Statistical Ecology). The authors are very grateful to the editor-in-chief, an associate editor, and two anonymous referees, whose comments have greatly improved the clarity of presentation of this material.

#### ORCID

Svetlana V. Tishkovskaya D https://orcid.org/0000-0003-3087-6380

## REFERENCES

Barnes, T. G. (2000). Landscape ecology and ecosystems management (Vol. 9, pp. 1–9). University of Kentucky College of Agriculture, Lexington and Kentucky State University.

Blackwell, P. G. (1997). Random diffusion models for animal movement. Ecological Modelling, 100, 87–102.

Blackwell, P. G. (2003). Bayesian inference for Markov processes with diffusion and discrete components. Biometrika, 90(3), 613-627.

Blackwell, P. G., & Macdonald, D. W. (2000). Shapes and sizes of badger territories. Oikos, 89, 392-398.

Blackwell, P. G., & Møller, J. (2003). Bayesian analysis of deformed tessellation models. Advances in Applied Probability, 35(1), 4–26.

Blackwell, P. G., Niu, M., Lambert, M. S., & LaPoint, S. D. (2016). Exact Bayesian inference for animal movement in continuous time. *Methods in Ecology and Evolution*, 7, 184–195. https://doi.org/10.1111/2041-210X.12460

# 16 of 17 | WILEY-

- Brillinger, D. R., Preisler, H. K., Ager, A. A., & Kie, J. G. (2001). *The use of potential functions in modeling animal movement Data Analysis from Statistical Foundations* (pp. 369–386). Nova science publishers, Inc.
- Cagnacci, F., Boitani, L., Powell, R. A., & Boyce, M. S. (2010). Animal ecology meets GPS-based radiotelemetry: A perfect storm of opportunities and challenges. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 365, 2157–2162.
- Calenge, C., Dray, S., & Royer-Carenzi, M. (2009). The concept of animals' trajectories from a data analysis perspective. *Ecological Informatics*, 4(1), 34–41.
- Christ, A., ver Hoef, J., & Zimmerman, D. (2008). An animal movement model incorporating home range and habitat selection. *Environmental and Ecological Statistics*, *15*, 27–38.
- Dunn, J. E., & Gipson, P. S. (1977). Analysis of radio telemetry data in studies of home range. Biometrics, 33, 85-101.
- Eisenhauer, E., & Hanks, E. (2020). A lattice and random intermediate point sampling design for animal movement. *Environmetrics*, *31*. https://doi.org/10.1002/env.2618
- Frair, J., Fieberg, J., Hebblewhite, M., Cagnacci, F., DeCesare, N., & Pedrotti, L. (2010). Resolving issues of imprecise and habitat-biased locations in ecological analyses using GPS telemetry data. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 365, 2187–2200.
- Gelman, A., Carlin, B. P., Stern, H. S., & Rubin, D. B. (2004). Bayesian data analysis (2nd ed.). Chapman & Hall.
- Gloaguen, P., Mahevas, S., Rivot, E., Woillez, M., Guitton, J., Vermard, Y., & Etienne, M. P. (2015). An autoregressive model to describe fishing vessel movement and activity. *Environmetrics*, 26(1), 17–28. https://doi.org/10.1002/env.2319
- Gurnell, J. (1984). Home range, territoriality, caching behaviour and food supply of the red squirrel (*Tamiasciurus hudsonicus fremonti*) in a subalpine lodgepole pine forest. *Animal Behaviour*, 32(4), 1119–1131.
- Harris, K. J. (2007). Statistical modelling and inference for radio-tracking (PhD thesis). The University of Sheffield.
- Harris, K. J., & Blackwell, P. G. (2013). Flexible continuous-time modelling for heterogeneous animal movement. *Ecological Modelling*, 255, 29–37.
- Hooten, M. B., Johnson, D. S., McClintock, B. T., & Morales, J. M. (2017). Animal movement: Statistical models for telemetry data. CRC Press.
- Johnson, D. S., London, J. M., Lea, M.-A., & Durban, J. W. (2008). Continuous-time correlated random walk model for animal telemetry data. *Ecology*, *89*(5), 1208–1215.
- Johnson, D. S., Thomas, D. L., Ver Hoef, J. M., & Christ, A. (2008). A general framework for the analysis of animal resource selection from telemetry data. *Biometrics*, 64(3), 968–976.
- Linn, I., Perrin, M. R., & Hiscocks, K. (2007). Use of space by the four-toed elephant-shrew *Petrodromus tetradactylus* (Macroscelidae) in Kwazulu-Natal (South Africa). *Mammalia*, 71(1/2), 30–39.
- McClintock, B. T. (2017). Incorporating telemetry error into hidden Markov models of animal movement using multiple imputation. *Journal of Agricultural Biological and Environmental Statistics*, 22(3), 249–269. https://doi.org/10.1007/s13253-017-0285-6
- Michelot, T., & Blackwell, P. G. (2019). State-switching continuous-time correlated random walks. *Methods in Ecology and Evolution*, 10, 637–649.
- Milner, J. E., Blackwell, P. G., & Niu, M. (2021). Modelling and inference for the movement of interacting animals. *Methods in Ecology and Evolution*.
- Niu, M., Blackwell, P. G., & Skarin, A. (2016). Modelling interdependent animal movement in continuous time. *Biometrics*, 72(2), 315–324. https://doi.org/10.1111/biom.12454
- Patterson, T. A., Parton, A., Langrock, R., Blackwell, P. G., Thomas, L., & King, R. E. (2016). Statistical modelling of individual animal movement: An overview of key methods and a discussion of practical challenges. *Advances in Statistical Analysis*, *101*, 399–438.
- Pope, C. A., Gosling, J. P., Barber, S., Johnson, J. S., Yamaguchi, T., Feingold, G., & Blackwell, P. G. (2019). Gaussian process modeling of heterogeneity and discontinuities using Voronoi tessellations. *Technometrics*, 63, 53–63. https://doi.org/10.1080/00401706.2019.1692696
- Preisler, H. K., Ager, A. A., Johnson, B. K., & Kie, J. G. (2004). Modeling animal movements using stochastic differential equations. *Environmetrics*, 15, 643–657. https://doi.org/10.1002/env.636
- Russell, J. C., Hanks, E. M., Haran, M., & Hughes, D. (2018). A spatially varying stochastic differential equation model for animal movement. Annals of Applied Statistics, 12(2), 1312–1331. https://doi.org/10.1214/17-AOAS1113
- Russell, R. E., Royle, J. A., Desimone, R., Schwartz, M. K., Edwards, V. L., Pilgrim, K. P., & Mckelvey, K. S. (2012). Estimating abundance of mountain lions from unstructured spatial sampling. *The Journal of Wildlife Management*, 76, 1551–1561. https://doi.org/10.1002/jwmg.412
- Scharf, H., Hooten, M. B., & Johnson, D. S. (2017). Imputation approaches for animal movement modeling. *Journal of Agricultural Biological and Environmental Statistics*, 22(3), 335–352. https://doi.org/10.1007/s13253-017-0294-5
- Scharf, H. R., Hooten, M. B., Wilson, R. R., Durner, G. M., & Atwood, T. C. (2019). Accounting for phenology in the analysis of animal movement. *Biometrics*, 75, 810–820.
- Spiegelhalter, D. J., Best, N. G., Carlin, B. P., & van der Linde, A. (2002). Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society, Series B*, 64, 583–640.
- Tanner, M., & Wong, W. H. (1987). The calculation of posterior distribution by data augmentation. *Journal of the American Statistical* Association, 82(398), 528-540.

- Wang, Y., Blackwell, P. G., Merkle, J. A., & Potts, J. R. (2019). Continuous time resource selection analysis for moving animals. *Methods in Ecology and Evolution*, 10, 1664–1678.
- Wilson, R. R., Hooten, M. B., Strobel, B. N., & Shivik, J. A. (2010). Accounting for individuals, uncertainty, and multiscale clustering in Core area estimation. *Journal of Wildlife Management*, 74(6), 1343–1352. https://doi.org/10.2193/2009-438

Worton, B. J. (1995). Modelling radio-tracking data. Environmental and Ecological Statistics, 2, 15-23.

#### SUPPORTING INFORMATION

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

How to cite this article: Tishkovskaya SV, Blackwell PG. Bayesian estimation of heterogeneous environments from animal movement data. *Environmetrics*. 2021;e2679. https://doi.org/10.1002/env.2679