

Efficient Bayesian analysis of occupancy models with logit link functions

Allan E. Clark^{1,2}  | Res Altwegg^{1,2}

¹Department of Statistical Sciences, University of Cape Town, Cape Town, South Africa

²Center for Statistics in Ecology, Environment and Conservation (SEEC), University of Cape Town, Rondebosch, South Africa

Correspondence

Allan E. Clark, Department of Statistical Sciences, University of Cape Town, Cape Town, South Africa.

Email: allan.clark@uct.ac.za

Funding information

National Research Foundation of South Africa, Grant/Award Number: 81685 and 99385

Abstract

Occupancy models (Ecology, 2002; 83: 2248) were developed to infer the probability that a species under investigation occupies a site. Bayesian analysis of these models can be undertaken using statistical packages such as *WinBUGS*, *OpenBUGS*, *JAGS*, and more recently *Stan*, however, since these packages were not developed specifically to fit occupancy models, one often experiences long run times when undertaking an analysis. Bayesian spatial single-season occupancy models can also be fit using the R package *stocc*. The approach assumes that the detection and occupancy regression effects are modeled using probit link functions. The use of the logistic link function, however, is algebraically more tractable and allows one to easily interpret the coefficient effects of an estimated model by using odds ratios, which is not easily done for a probit link function for models that do not include spatial random effects. We develop a Gibbs sampler to obtain posterior samples from the posterior distribution of the parameters of various occupancy models (nonspatial and spatial) when logit link functions are used to model the regression effects of the detection and occupancy processes. We apply our methods to data extracted from the 2nd Southern African Bird Atlas Project to produce a species distribution map of the Cape weaver (*Ploceus capensis*) and helmeted guineafowl (*Numida meleagris*) for South Africa. We found that the Gibbs sampling algorithm developed produces posterior samples that are identical to those obtained when using *JAGS* and *Stan* and that in certain cases the posterior chains mix much faster than those obtained when using *JAGS*, *stocc*, and *Stan*. Our algorithms are implemented in the R package, *Rcppocc*. The software is freely available and stored on GitHub (<https://github.com/AllanClark/Rcppocc>).

KEYWORDS

Bayesian spatial occupancy model, imperfect detection, occupancy model, *Rcppocc*, restricted spatial regression

1 | INTRODUCTION

Occupancy models are an important statistical technique that was developed to make use of detection/nondetection data to infer the probability that a species under investigation occupies a site. When

an occupancy study is undertaken, n_s sites are visited a number of times to estimate the occupancy probability (ψ) and conditional detection probability (p) of a species associated with each site in a region. The method can be viewed as an extension of logistic regression and allows one to estimate the occupancy probability at

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.

© 2019 The Authors. *Ecology and Evolution* published by John Wiley & Sons Ltd.

sites where none of the species being investigated have been detected. The model is formulated hierarchically, using Bernoulli random variables to specify the occupancy and detection processes, respectively, which can be modeled using site-specific and survey-specific explanatory variables, respectively (MacKenzie et al., 2002). Johnson, Conn, Hooten, Ray, and Pond (2013) note that occupancy models produce “unbiased inference when occupancy observations at nearby units are conditionally independent given any available covariates” but stress that “spatial autocorrelation may lead to biases and overestimated precision” of regression effects. This observation has led to the development of various models to account for spatial autocorrelation in ecological survey data (Aing, Halls, Oken, Dobrow, & Fieberg, 2011; Gardner, Lawler, Ver Hoef, Magoun, & Kellie, 2010; Hoeting, Leecaster, & Bowden, 2000; Hooten, Larsen, & Wikle, 2003) and have extensively been used to guide environmental monitoring and assessment programs globally.

A number of methods have been used to fit occupancy models to data. These include maximum likelihood (MacKenzie et al., 2002); penalized maximum likelihood (Hutchinson, Valente, Emerson, Betts, & Dietterich, 2015; Moreno & Lele, 2010), Bayesian methods that employ *WinBUGS*, *OpenBUGS*, *JAGS*, or *Stan* as well as approximate methods such as those developed by Clark, Altwegg, and Ormerod (2016). Recently Dorazio and Rodriguez (2012) and Johnson et al. (2013) developed Gibbs algorithms to obtain posterior samples for the parameters of a nonspatial and spatial single-season occupancy (SSO) model, respectively. Both approaches assume that detection and occupancy processes are modeled using probit link functions, which enables the use of data augmentation (Tanner & Wong, 1987) to obtain closed form expressions of the conditional posterior distributions of the parameters of the occupancy model.

Given that the probit and logistic functions are very similar and only differ in respect of the tails of the functions, analysis undertaken using either of the functions should produce similar occupancy and conditional detection probabilities (Dorazio & Rodriguez, 2012). However, the use of the logistic link function is algebraically more tractable and allows one to easily interpret the coefficient effects of an estimated model by using odds ratios, which is not easily done for a probit link function. This observation is particularly true for the nonspatial SSO model since no spatial random effects are included in this model; however, when spatial random effects are included in the model, the interpretation of the regression effects can be difficult (Boehm, Reich, & Bandyopadhyay, 2013).

The paper commences with a brief discussion of the link between logistic regression and occupancy models. Thereafter, we discuss the formulation of various popular Bayesian spatial occupancy models and develop a Gibbs sampling algorithm for a particular spatial occupancy model when the regression effects of the occupancy and detection processes are modeled using logit link functions. Before concluding, we analyze two detection/nondetection data sets of South African bird species to illustrate the methods developed in the paper. An R package (*Rcppocc*) has been developed to fit SSO models using Gibbs sampling which can be obtained at: <https://github.com/AllanClark/Rcppocc>.

2 | MATERIAL AND METHODS

2.1 | Logistic regression and occupancy models

Assume that n_s sites are surveyed a number of times and detection/nondetection data are collected at all sites. Denote the observed data as a ragged matrix $\mathbf{y} = [y_{ij}]$ where $y_{ij} = 1$ if the species under investigation has been observed at site i during survey j and $y_{ij} = 0$ otherwise. Let the vector \mathbf{z} represents the true species occupancy at the sites considered such that $z_i = 1$ if the species occupies site i and $z_i = 0$ if it does not occupy site i . The SSO model can be represented using the following hierarchical model, $z_i | \psi_i \sim \text{Bernoulli}(\psi_i)$, $y_{ij} | z_i, p_{ij} \sim \text{Bernoulli}(z_i p_{ij})$ for all sites $i = 1, \dots, n_s$; for all surveys $j = 1, \dots, V_i$ (Royle & Dorazio, 2008). The variable ψ_i denotes the probability occupancy probability at site i , while $p_{ij} = \Pr(y_{ij} = 1 | z_i = 1)$ denotes the conditional probability of detecting the species during the j th survey of site i given that the species is present at site i . In what follows we assume that the conditional detection and occupancy regression effects (α and β) are modeled using logit link functions such that $\text{logit}(\psi_i) = \mathbf{x}_i^T \beta$ and $\text{logit}(p_{ij}) = \mathbf{w}_{ij}^T \alpha$, where \mathbf{x}_i^T and \mathbf{w}_{ij}^T are row vectors in design matrices, \mathbf{X} (occupancy) and \mathbf{W} (detection), respectively (as defined in Clark et al., 2016).

The joint posterior distribution of the parameters of the model is

$$[\mathbf{z}, \alpha, \beta | \mathbf{y}] \propto \pi(\alpha) \pi(\beta) \left(\prod_{i=1}^{n_s} \psi_i^{z_i} (1 - \psi_i)^{1 - z_i} \right) \prod_{i=1}^{n_s} \prod_{j=1}^{V_i} p_{ij}^{y_{ij}} (1 - p_{ij})^{1 - y_{ij}}$$

where $\pi(\alpha)$ and $\pi(\beta)$ are the prior distributions of α and β , respectively. A directed acyclic graph of the above problem is displayed in Figure 1 below.

A Gibbs sampling algorithm for the parameters of this model requires sampling from $[\beta | \mathbf{z}]$, $[\alpha | \mathbf{z}, \mathbf{y}]$, and $[z_i = 1 | \alpha, \beta, \mathbf{y}]$ for all sites where

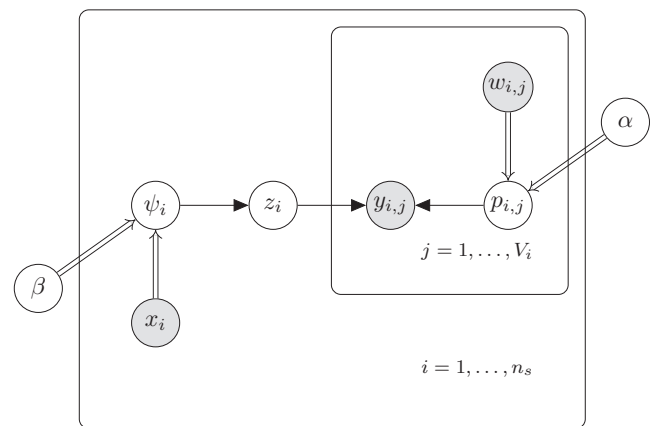


FIGURE 1 A directed acyclic graph illustrating the dependencies between the parameters and observed data for an SSO model. Shaded nodes represent observed data while all latent parameters are represented using unshaded nodes. Deterministic relationships are represented using double arrows while all stochastic relationships are represented using a single arrow

the species has not been observed. The first two conditional distributions have the following form,

$$[\beta|\mathbf{z}] \propto \pi(\beta) \prod_{i=1}^{n_s} \psi_i^{z_i} (1-\psi_i)^{1-z_i} \text{ and } \quad (1)$$

$$[\alpha|\mathbf{z}, \mathbf{y}] \propto \pi(\alpha) \prod_{i=1}^{n_s} \prod_{j \in \{i: z_i=1\}} p_{ij}^{y_{ij}} (1-p_{ij})^{1-y_{ij}}. \quad (2)$$

Notice that Equations 1 and 2 are of the same form as the posterior distributions of the regression effects of a logistic regression model and therefore we adapt a Gibbs sampling scheme for logistic regression models to address the problem of obtaining posterior samples for the parameters of an occupancy model.

In a logistic regression context, Polson, Scott, and Windle (2013) show that posterior samples of the regression effects can be obtained by sampling from the conditional distributions of Pólya-Gamma random variables and multivariate Gaussian distributions in turn. Their method is similar to that of Albert and Chib (1993) who developed a Gibbs algorithm to undertake probit regression, the only difference being that the sampling from truncated Gaussian distributions is replaced by sampling from Pólya-Gamma distributions. The sampling methods developed by Polson et al. (2013) are exact since their Pólya-Gamma sampling method is uniformly ergodic and converges to the correct posterior distribution (Choi, & Hobert, 2013).

For the SSO model, the conditional posterior distributions of $\beta|\omega_\beta, \mathbf{y}$ are derived by introducing Pólya-Gamma latent variables, ω_β , and noting that the contribution of the i^{th} observation to a Bernoulli likelihood can be re-expressed as

$$\begin{aligned} \psi_i^{z_i} (1-\psi_i)^{1-z_i} &= \frac{\exp(\mathbf{x}_i^T \beta)^{z_i}}{1 + \exp(\mathbf{x}_i^T \beta)} \\ &= \exp(\kappa_i \mathbf{x}_i^T \beta) \int \exp\left(-\frac{\omega_{i,\beta}}{2} (\mathbf{x}_i^T \beta)^2\right) p(\omega_{i,\beta} | 1, 0) d\omega_{i,\beta} \end{aligned}$$

where $\kappa_i = z_i - 0.5$ and $p(\omega_{i,\beta} | 1, 0)$ is the probability density function of a Pólya-Gamma distribution with parameters 1 and 0 (Polson et al., 2013). The conditional posterior distribution of α is derived by using the same manipulation of the Bernoulli likelihood.

In the Supporting information (Appendix S1), we discuss the existing Gibbs algorithms used for undertaking logistic regression and demonstrate the use of the Pólya-Gamma (PG) method by developing two Gibbs sampling algorithms for the parameters of SSO models. In Table 1, we summarize the Gibbs algorithms for an SSO model when using the PG method but provide the details regarding the algorithm in the Supporting Information (Appendix S2 and Appendix S3). We use the notation " $a \sim \text{PG}(b, c)$ " to indicate that the random variable a is a Pólya-Gamma random variable with parameters b and c . Take note that the algorithm is identical to that developed by Dorazio and Rodriguez (2012) except that the sampling from truncated Gaussian distributions is replaced by sampling from Pólya-Gamma distributions.

2.2 | Bayesian spatial SSO models

Spatial generalized linear mixed models (SGLMM) are an extension of the general linear model (Nelder & Wedderburn, 1972) that

TABLE 1 The Gibbs algorithm for undertaking a SSO model using the "PG" method (See the Supporting Information (Appendix S3) for the details pertaining to the parameter matrices of the conditional posterior distributions.)

-
- 1: Set starting values for α , β and \mathbf{z} .
 - 2: **for** (iterations = 1, ..., simulation runs) **do**
 - 3: **for** ($i = 1, \dots, n_s$) **do**
 - 4: - Generate $\omega_{i,\beta} \sim \text{PG}(1, \mathbf{x}_i^T \beta)$.
 - 5: **end for**
 - 6: - Generate $\beta \sim \mathcal{N}(\mu_\beta, \Sigma_\beta)$.
 - 7: **for** ($(ij) = \{z_i = 1\}$) **do**
 - 8: - Generate $\omega_{ij,\alpha} \sim \text{PG}(1, \mathbf{w}_{ij}^T \alpha)$.
 - 9: **end for**
 - 10: - Generate $\alpha \sim \mathcal{N}(\mu_\alpha, \Sigma_\alpha)$.
 - 11: - Generate $z_i \sim \text{Bin}(\tilde{\psi}_i)$.
 - 12: **end for**
-

allows the link function of the expected value of the random variable under investigation to be modeled as a function of a spatial random variable/s. The formulation was first developed by Besag, York, and Mollié (1991) and has been extensively used in areas such as agriculture (Besag & Higdon, 1999), biostatistics (Gelfand & Vounatsou, 2003; Waller & Gotway, 2004), ecology (Lichstein, Simons, Shiner, & Franzreb, 2002) and species distribution modelling (Drouilly, Clark, & O'Riain, 2018; Gelfand et al., 2005; Hooten et al., 2003; Latimer, Wu, Gelfand, & Silander, 2006).

The paper by Gelfand et al. (2005) led to the development of the R package *hSDM* (Vieilledent et al., 2014) in which the *hSDM.siteocc.iCAR* function can be used to fit a particular spatial occupancy model to detection/nondetection data. A region under investigation is subdivided into n_s grid cells each which are surveyed a number of occasions. The model is formulated using Bernoulli latent random variables $\mathbf{z} = (z_1, \dots, z_{n_s})^T$. Formally, $z_i | \psi_i \sim \text{Bernoulli}(\psi_i)$ with $\text{logit}(\psi_i) = \mathbf{x}_i^T \beta + \rho_i$, for all $i = 1, \dots, n$, where $\rho = (\rho_1, \dots, \rho_{n_s})^T$ is a multivariate Gaussian random vector with mean $\mathbf{0}$ and correlation matrix defined using the neighborhood structure of the grid cells. The observation process is specified as in the non-spatial model. The documentation of the function indicates that posterior samples of the parameters of the model are obtained using the C programming language and utilizes an adaptive Metropolis algorithm (Metropolis, Rosenbluth, Rosenbluth, Teller, & Teller, 1953; Robert & Casella, 2013).

Johnson et al. (2013) develop two spatial occupancy models. They assume that probit link functions are used to model both the occupancy and detection processes and thereby rely on data augmentation to develop a Gibbs sampling algorithm to sample from the posterior distribution of the parameter of the models. For the probit case, the occupancy probability of a particular grid cell (for the standard occupancy model) is calculated as $\Phi(\mathbf{x}_i^T \beta) = \Pr(z_i = 1)$. In a Bayesian context, such a probit model is formulated by defining a latent Gaussian

random variable, \tilde{z}_i with mean 0 and variance equal to 1 such that $\Pr(z_i = 1) = \Pr(\tilde{z}_i > 0)$. In the first of their models, they allow \tilde{z}_i to be spatially correlated such that $\tilde{z}_i = \mathbf{x}_i^T \boldsymbol{\beta} + \eta_i + \epsilon_i$. $\boldsymbol{\eta} = (\eta_1, \dots, \eta_{n_s})^T$ is defined as $\boldsymbol{\rho}$ above while $\epsilon_i \sim \mathcal{N}(0,1)$, for all $i = 1, \dots, n_s$.

Often the spatial random effects and fixed effects of a model are collinear when spatially varying covariates are included as fixed effects (Hanks, Schliep, Hooten, & Hoeting, 2015; Hodges & Reich, 2010; Hughes & Haran, 2013; Reich, Hodges, & Zadnik, 2006). The suggested solution to this problem was to include spatial random effects in the model specification that are orthogonal to the fixed effects and is known as *restricted spatial regression* (RSR). The second spatial model developed by Johnson et al. (2013) uses this method and redefines \tilde{z}_i as $\tilde{z}_i = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{k}_i^T \boldsymbol{\theta} + \epsilon_i$, where \mathbf{k}_i^T is a row vector of the design matrix \mathbf{K} . The spatial random effects are modeled as

$$\theta | \tau \sim \mathcal{N}\left(\mathbf{0}_r, \frac{1}{\tau} (\mathbf{K}^T \mathbf{Q} \mathbf{K})^{-1}\right) = \mathcal{N}\left(\mathbf{0}_r, \frac{1}{\tau} \mathbf{M}\right),$$

$$\tau \sim \mathcal{G}(i_1, i_2) \text{ and } \epsilon \sim \mathcal{N}(\mathbf{0}_n, \mathbf{I}_n).$$

\mathbf{Q} is a $n \times n$ ICAR precision matrix (Besag & Kooperberg, 1995) obtained using surveyed and unsurveyed locations, τ is a spatial precision parameter and i_1 and i_2 are known constants. Kelsall and Wakefield (1999) have suggested setting these parameters to 0.5 and 0.005, respectively, such that the prior mean of τ is 1,000. The matrix \mathbf{K} consists of the first r ($r \ll n$) eigenvectors of $\boldsymbol{\Omega} = \mathbf{n} \mathbf{R} \mathbf{A} \mathbf{R}^T / \mathbf{1}^T \mathbf{A} \mathbf{1}$ where $\mathbf{R} = \mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ and \mathbf{A} is an association matrix with (ij) th entry $A_{ij} = 1$ if sites i and j are neighbours and zero otherwise.

In our formulation of the spatial occupancy model, we model the occupancy probabilities at all grid cells as

$$\text{logit}(\psi_i) = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{k}_i^T \boldsymbol{\theta}, \quad (3)$$

where \mathbf{K} and $\boldsymbol{\theta}$ are defined above. We make use of Pólya-Gamma random variables to obtain the conditional distributions of the parameters of the above spatial occupancy model. The conditional distributions are very similar to those obtained for the SSO model although here we require the conditional posterior distribution of additional parameters ($\boldsymbol{\theta}$ and τ). A directed acyclic graph of the spatial SSO model is displayed in Figure 2 below while in Table 2, we summarize the Gibbs algorithms for a spatial SSO model which employs Equation 3 when using the PG method. The details regarding the algorithm can be found in the Supporting information (Appendix S4).

2.3 | Applications

To demonstrate our methods, we used detection/nondetection data extracted from the 2nd Southern African Bird Atlas Project (SABAP2) database to produce a species distribution map of the Cape weaver (*Ploceus capensis*) and helmeted guineafowl (*Numida meleagris*) for South Africa. SABAP2 divides Southern Africa into a continuous grid of $5' \times 5'$ and relies on citizen scientists to collect checklists of bird species for each grid cell. Birders are requested

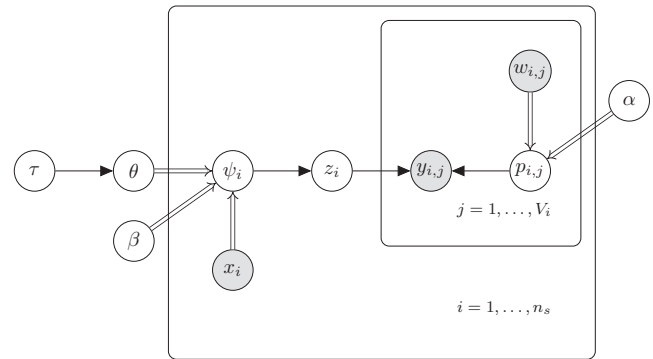


FIGURE 2 A directed acyclic graph illustrating the dependencies between the parameters and observed data for a spatial SSO model. Shaded nodes represents observed data while all latent parameters are represented using unshaded nodes. Deterministic relationships are represented using double arrows while all stochastic relationships are represented using a single arrow

TABLE 2 The Gibbs algorithm for undertaking a spatial SSO model (See the Supporting Information (Appendix S4) for the details pertaining to the parameter matrices of some of the conditional posterior distributions.)

- 1: Set starting values for α , β , $\boldsymbol{\theta}$ and \mathbf{z} .
- 2: **for** (iterations = 1, ..., simulation runs) **do**
- 3: **for** ($i = 1, \dots, n_s$) **do**
- 4: - Generate $\omega_{i,\beta} \sim \text{PG}(1, \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{k}_i^T \boldsymbol{\theta})$.
- 5: **end for**
- 6: - Generate $\beta \sim \mathcal{N}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta)$.
- 7: - Generate $\boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\mu}_\theta, \boldsymbol{\Sigma}_\theta)$.
- 8: - Generate $\tau \sim \mathcal{G}\left(\frac{r}{2} + i_1, \frac{\boldsymbol{\theta}^T \mathbf{M}^{-1} \boldsymbol{\theta}}{2} + i_2\right)$.
- 9: **for** ($(ij) = \{z_i = 1\}$) **do**
- 10: - Generate $\omega_{ij,\alpha} \sim \text{PG}(1, w_{ij}^T \alpha)$.
- 11: **end for**
- 12: - Generate $\alpha \sim \mathcal{N}(\boldsymbol{\mu}_\alpha, \boldsymbol{\Sigma}_\alpha)$.
- 13: - Generate $z_i \sim \text{Bin}(\tilde{\psi}_i)$.
- 14: **end for**

to spend at least 2 hr on each checklist in which they undertake *intense birding* and record all species they observe and the order in which they are observed. For this analysis, we aggregated the data to quarter-degree grid cells. We used data that span South Africa and contained a minimum of three and a maximum of fifty surveys during 2016 (January–December) in the analysis. Covariate information at unsurveyed locations was included in the analyses to obtain occupancy estimates that span South Africa. All covariates were centered and standardized.

TABLE 3 Posterior run times for the Bayesian spatial occupancy models as well as the ESR (per minute) for α , β , and τ

Species	Method	Time (min)	α_0	α_1	β_0	β_1	β_2	τ
Cape weaver	<i>stocc</i> (RSR)	27.00	552.71	725.59	12.03	7.83	22.37	11.35
	<i>stocc</i> (ICAR)	136.76	104.41	142.57	0.25	0.10	0.15	0.14
	JAGS	243.55	102.62	131.42	2.69	1.71	5.09	1.97
	<i>Stan</i>	187.06	275.75	331.44	53.48	34.42	83.74	19.35
	<i>Rcppocc</i>	19.88	1682.15	1804.42	85.32	53.40	141.82	116.19
Helmeted Guinea fowl	<i>stocc</i> (RSR)	27.08	619.72	563.24	11.63	43.61	48.34	15.61
	<i>stocc</i> (ICAR)	165.23	65.92	97.15	0.04	0.16	0.09	0.05
	JAGS	254.49	97.86	125.35	3.08	10.84	12.06	2.66
	<i>Stan</i>	150.55	595.77	617.84	49.21	186.23	286.92	24.44
	<i>Rcppocc</i>	19.9	1888.11	1493.07	86.16	314.04	364.47	106.06

In our analysis, we fitted a nonspatial and spatial SSO model with one detection covariate and two occupancy covariates. The detection covariate used was the number of species observed by the birder (denoted as *nssp*), while the occupancy covariates were functions of seven climate variables. It is assumed that the more species the birder observes while birding, the more likely they are to observe the particular species being analyzed such that a positive detection regression effect is expected. The aim of the analysis is not to obtain the best occupancy model for the particular data sets but rather to highlight the use of the developed Gibbs sampling algorithm for fitting RSR occupancy models to the data using different software programs and different sampling methods. We specifically consider the Gibbs sampling algorithm by Johnson et al. (2013) (probit link functions), our Gibbs algorithm (logit link functions), JAGS, and *Stan* (which uses a *no-U-turn* Hamiltonian Monte Carlo sampler (Hoffman & Gelman, 2014) to sample from the parameters of a posterior distribution).

The climate variables (Figure S5 in the Supporting Information Appendix S6) all form part of a data set used by Huntley et al. (2006) to model bird distributions in Southern Africa. The variables included two measures of annual temperature that related to thermal sums above 0 and 5 degree centigrade; two measures related to the mean temperature of the coldest and warmest month, respectively; the ratio of potential to realized evapotranspiration as well as two measures that relates to the intensity of the dry and wet season, respectively. The climate variables are highly correlated with two of the variables having variance inflation factors in excess of 3,000 (Tables S2 and S3 in the Supporting Information Appendix S5). Because of this fact, it was decided to extract two principal components from the design matrix that consisted of the centered and standardized climate variables. These principal components explain 90% of the variation in the design matrix (Table S4 in the Supporting Information Appendix S5) and can tentatively be interpreted as a temperature related factor and a climate intensity factor, respectively.

We follow Hughes and Haran (2013) and retain 10% of the eigenvalues ($\lambda_i, i = 1, \dots, n$) of Ω . In a similar context, Johnson et al. (2013) suggest selecting a RSR model with $\lambda_i \geq 0.5$ which suggests including at most 237 eigenvectors into the spatial portion of the model.

Experimentation with different values of r between 150 and 230 demonstrated no significant difference to the results we report here.

The following prior distributions were used for the parameters of the spatial SSO model: $\alpha \sim \mathcal{N}(\mathbf{0}, 1000I_2)$, $\beta \sim \mathcal{N}(\mathbf{0}, 1000I_3)$, and $\tau \sim \mathcal{G}(0.5, 0.005)$. The prior specification for τ places more weight on large values of τ indicating that very little prior weight is placed on the spatial random effects of the model. Broms (2013) performed a simulation study and found that the RSR model results are not sensitive to the prior specification of τ and thus we have not done any analysis to test the sensitivity of our results to the prior specification of τ . All MCMC sampling was undertaken using the R packages, *stocc*, *jagsUI* (Kellner, 2014) in combination with JAGS 4.2.0 (Plummer, 2003), *rstan* in combination with *Stan* 2.17.3 as well as the authors' code.¹ All calculations were performed on a Windows 10 Pro desktop computer which had an Intel(R) Core(TM) i7-6900 processor with 64 GB of RAM. One chain of 70,000 iterations was run. The first 20,000 samples were discarded as burn-in samples, while the remaining samples were retained. Experimentation and an examination of the Geweke convergence diagnostic statistics (Geweke, 1992) and trace plots obtained by running three parallel chains using *Rcppocc* displayed that the MCMC chains converged using these numbers of iterations. The posterior samples were not thinned (Link & Eaton, 2012).

3 | RESULTS

From the analysis of both data sets, we observe that the Gibbs algorithm developed for the spatial occupancy model produces identical posterior distributions to those obtained when using JAGS and *Stan* (Figures A1 and A2 in Appendix 1). In both data sets, the posterior samples of the detection regression effects exhibit good mixing where the lagged sample autocorrelations of the posterior samples approach zero within 5 lags. The posterior samples of the occupancy regression effects as well as the precision of the spatial random effect (τ) exhibit slower mixing when using *stocc*, JAGS, and *Rcppocc* (denoted as "PG" in Figures A3 and A4 in Appendix 2), while *Stan* produced a posterior chain that mixed well. We observe that *stocc*

TABLE 4 Posterior summaries of the parameters of the Bayesian nonspatial and spatial occupancy models (posterior mean, Monte Carlo standard error, standard deviation, 2.5% and 97.5% quantiles)

Type	Species	Parameter	Mean	MCSE	SD	2.5%	97.5%
Nonspatial	Cape weaver	α_0	-0.32	0.0002	0.03	-0.38	-0.26
		α_1	0.56	0.0002	0.03	0.49	0.62
		β_0	-0.49	0.0006	0.10	-0.68	-0.30
		β_1	-0.71	0.0005	0.06	-0.84	-0.59
		β_2	-0.24	0.0004	0.06	-0.36	-0.12
	Helmeted guineafowl	α_0	-0.10	0.0001	0.02	-0.15	-0.06
		α_1	0.78	0.0002	0.03	0.72	0.84
		β_0	-0.35	0.0006	0.06	-0.48	-0.22
		β_1	-0.36	0.0005	0.09	-0.54	-0.20
		β_2	-0.39	0.0008	0.08	-0.56	-0.24
Spatial	Cape weaver	α_0	-0.33	0.0001	0.03	-0.39	-0.27
		α_1	0.58	0.0001	0.03	0.52	0.64
		β_0	-1.54	0.0030	0.30	-2.17	-1.00
		β_1	-1.51	0.0027	0.20	-1.96	-1.16
		β_2	-0.55	0.0012	0.15	-0.86	-0.27
		τ	0.02	0.0001	0.01	0.01	0.03
	Helmeted guineafowl	α_0	-0.10	0.0001	0.02	-0.15	-0.05
		α_1	0.80	0.0001	0.03	0.74	0.85
		β_0	1.85	0.0029	0.25	1.40	2.40
		β_1	-0.36	0.0005	0.09	-0.54	-0.20
		β_2	-0.36	0.0004	0.10	-0.56	-0.17
		τ	0.03	0.0002	0.01	0.01	0.05

produced posterior samples for α , β , and τ that had the largest levels of autocorrelation among all of the methods considered (when fitting the RSR model).

Table 3 tabulates the run times (in minutes) and effective sampling rate (ESR^2 = the effective sample size per unit run time) of α , β , and τ for each of the sampling algorithms used to analyze the two data sets. For completeness sake, we also include the statistics related to the ICAR model when using *stocc*. We observe that *stocc* and *Rcppocc* had faster run times than *JAGS* and *Stan*. *Rcppocc* had the fastest running times and completed the 70,000 MCMC iterations approximately 12 times faster than *JAGS* and between 7 and 10 times faster than *Stan*. *Rcppocc* has the largest ESR of all of the algorithms considered and produced ESR values which ranged between 1.5 and 6 times larger than those obtained by *Stan*; 3–11 times larger than those obtained by *stocc* and 11–60 times larger than those obtained by *JAGS*. The ICAR models took approximately 8 times longer to run than the RSR model when using *stocc* and resulted in significantly larger levels of autocorrelation within the α , β , and τ chains.

The posterior summaries for some of the parameters of the nonspatial and spatial model are displayed in Table 4. The fixed regression effects of all of the parameters (for both data sets) were statistically different from zero since none of the 95% highest

density credibility interval of the parameters contained zero. In all cases, the regression effect for *nspp* was positive (as expected), while the regression effects of the occupancy effects were negative. The detection regression effects for both model types (for the respective species) were identical. The regression effects for the occupancy process for the Cape weaver were significantly different for the two model types, while the same regression effects for the helmeted guineafowl were identical for both model types (except for the intercept). The posterior distribution of the spatial standard deviation parameter ($\sigma = 1/\sqrt{\tau}$) indicates that the spatial process does significantly contribute to the variability of the occupancy process across South Africa. The 95% posterior highest density credibility interval for σ is [5.56, 10.59] and [4.26, 8.42] for the Cape weaver and the helmeted guineafowl data sets, respectively.

Figure 3a,c displays the estimated occupancy probabilities ($\Pr(z_i = 1|.)$) across South Africa estimated using *Rcppocc* for the Cape weaver and helmeted guineafowl data set, respectively. The figures illustrate that there is a high probability that the Cape weaver occupies coastal regions throughout South Africa and low occupancy probability (close to zero) in the interior areas of South Africa. In contrast, the helmeted guineafowl has very high occupancy probabilities in most regions of South Africa except for the North West regions of South Africa. Figure 3b,d displays the difference between the estimated occupancy probabilities

obtained when using *Rcppocc* and *stocc*, respectively ("*Rcppocc-stocc*"). The figures illustrate that we obtain similar estimates of the mean occupancy probabilities when using either estimation method with small discrepancies at the majority of the grid cells across South Africa.

4 | DISCUSSION AND CONCLUSIONS

Through several studies, Bayesian methods have been developed to undertake occupancy models. They, however, either use probit

link functions to model the detection and occupancy processes of the model; use general Bayesian analysis software such as *JAGS*, *WinBUGS*, *OpenBUGS*, or *Stan* to undertake their analysis or make use of the Metropolis–Hastings algorithm to sample from the parameters of the model. We develop a Gibbs sampling algorithm to obtain posterior samples of the parameters of a restricted spatial regression (RSR) occupancy model and demonstrate that the method has a larger expected sampling rate (ESR) and faster run times when compared to previous Bayesian methods used in the literature to date.

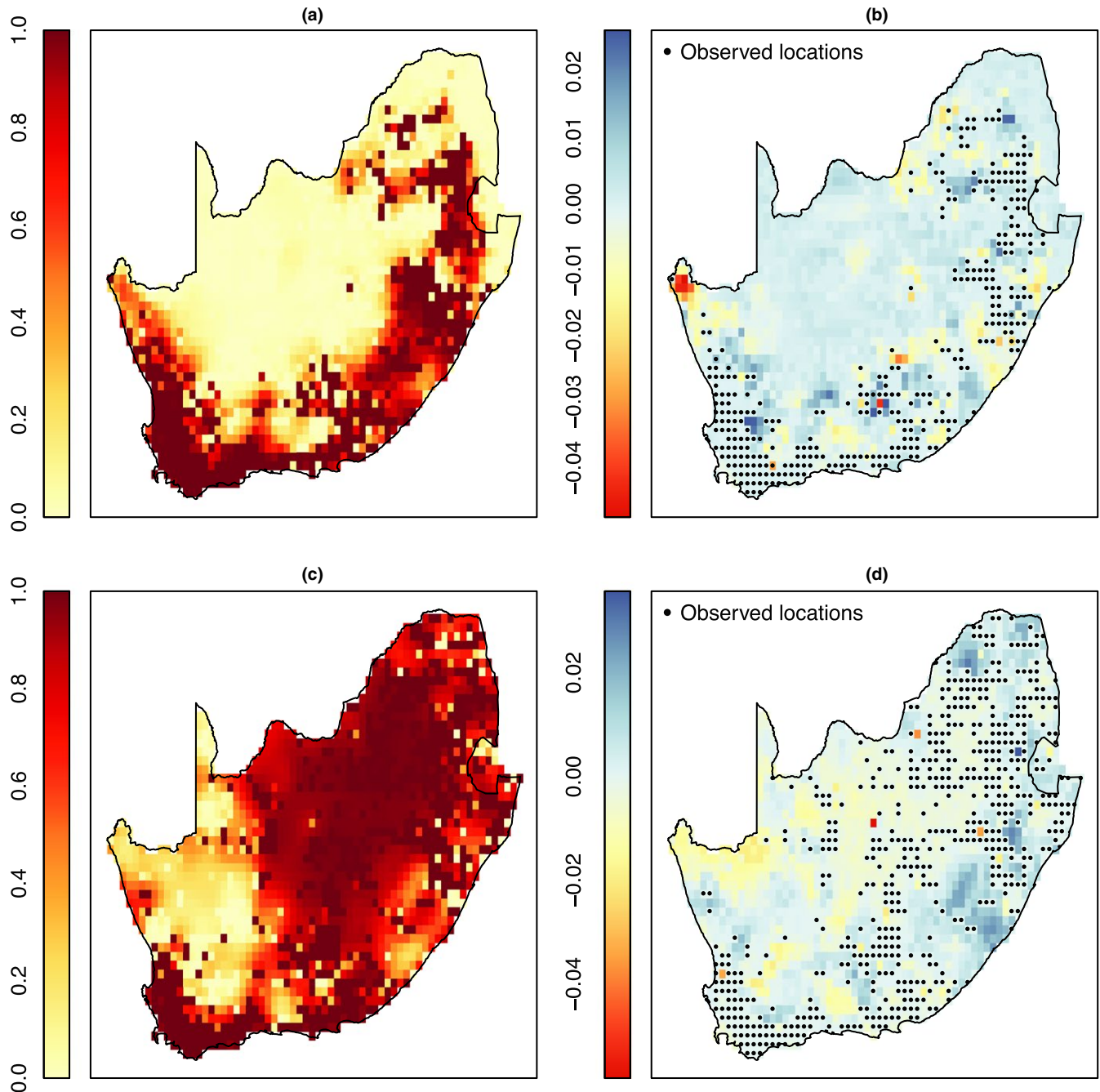


FIGURE 3 Estimated occupancy probability for the Cape weaver and helmeted guineafowl estimated using *Rcppocc* (a and c). The difference between the estimated occupancy probabilities obtained when using *Rcppocc* and *stocc* for the Cape weaver and helmeted guineafowl, respectively (b and d). The grid cells where the species have been detected at least once are displayed in (b) and (d)

Similar to Broms, Johnson, Altwegg, and Conquest (2014) and Johnson et al. (2013), we show that the ICAR model produced posterior samples with significantly larger autocorrelations than the RSR model when using *stocc*. As an example, the autocorrelations of the occupancy regression effects as well as the spatial precision parameter (of both data sets) had autocorrelations in excess of 0.7 at lag 500 indicating that the posterior chain of the model mixed poorly for those parameters of the ICAR model. Additionally, the run times of the ICAR model were approximately 5 times longer than the run times of the RSR model and thus we do not recommend its use when fitting a spatial occupancy model.

Based on the two data sets, we observed that the new algorithm not only ran faster (approximately 35%) than the Gibbs sampler implemented in *stocc*, it also generated expected sample size (ESS) statistics between 2 and 6 times larger than those obtained using *stocc*. The main reason for the time difference is that *stocc* has been coded using R, while *Rcppocc* uses Rcpp and RcppArmadillo to undertake all matrix computations. *Stan* uses compiled C++ code to implement the *no-U-turn* Hamiltonian Monte Carlo algorithm and generated ESS statistics between 2 and 7 times larger than those obtained using *Rcppocc*. In many applications, *Stan* has been shown to be much faster than JAGS although at present *Stan* has run times that are approximately 7–10 times slower than *Rcppocc* when fitting spatial occupancy models. The opportunity thus exists to develop suitable *Stan* (or NIMBLE) code that can fit spatial occupancy models in a shorter period of time.

ACKNOWLEDGMENTS

This research was partially supported by two South African National Research Foundation grants, namely, 99385 (Clark) and 81685 (Altwegg). The financial assistance of the NRF toward this research is hereby acknowledged. Opinions expressed and conclusions arrived at, are those of the author and are not necessarily to be attributed to the NRF. Allan Clark would also like to acknowledge the help of Andrew D. Crosby, a Postdoctoral Fellow at the Boreal Avian Modelling Project, Department of Biological Sciences, University of Alberta. He shared code (with Allan Clark) on how to fit the single-season occupancy model using *Stan* via email correspondence. The authors would also like to thank Prof Linda Haines (University of Cape Town) for reading the initial manuscript and providing helpful comments.

CONFLICT OF INTEREST

The authors have no conflict of interests to declare.

AUTHOR CONTRIBUTION

Below, Allan Ernest Clark is denoted as “AEC”, while Res Altwegg is denoted as “RA”. AEC and RA conceived and designed the paper. AEC analyzed the data. AEC wrote and, AEC and RA reviewed the paper. AEC designed and coded the software used in the analysis. AEC wrote computer code used to perform all analysis.

Notes

¹An R package has been developed to fit these models using MCMC. All code can be obtained from <https://github.com/AllanClark/Rcppocc>. Appendix S7 in the Supporting Information includes a worked example explaining how to run RSR models using *stocc*, *Stan*, and *Rcppocc*.

²ESR = the effective sample size per unit run time. The effective sample size for the *i*th parameter in the model is defined as $ESS_i = M/1 + 2 \sum_{j=1}^k \rho_i(j)$, where *M* is the number of retained samples, and $\rho_i(j)$ is the *j*th lagged autocorrelation of parameter *i* (Holmes & Held, 2006). We use the *coda* package (Plummer, Best, Cowles, & Vines, 2006) to estimate ESS_{*i*}.

ORCID

Allan E. Clark  <https://orcid.org/0000-0003-3472-0797>

REFERENCES

- Aing, C., Halls, S., Oken, K., Dobrow, R., & Fieberg, J. (2011). A bayesian hierarchical occupancy model for track surveys conducted in a series of linear, spatially correlated, sites. *Journal of Applied Ecology*, 48(6), 1508–1517. <https://doi.org/10.1111/j.1365-2664.2011.02037.x>
- Albert, J. H., & Chib, S. (1993). Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88(422), 669–679. <https://doi.org/10.1080/01621459.1993.10476321>
- Besag, J., & Higdon, D. (1999). Bayesian analysis of agricultural field experiments. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 61(4), 691–746. <https://doi.org/10.1111/1467-9868.00201>
- Besag, J., & Kooperberg, C. (1995). On conditional and intrinsic autoregressions. *Biometrika*, 82(4), 733–746.
- Besag, J., York, J., & Mollié, A. (1991). Bayesian image restoration, with two applications in spatial statistics. *Annals of the Institute of Statistical Mathematics*, 43(1), 1–20. <https://doi.org/10.1007/BF00116466>
- Boehm, L., Reich, B. J., & Bandyopadhyay, D. (2013). Bridging conditional and marginal inference for spatially referenced binary data. *Biometrics*, 69(2), 545–554. <https://doi.org/10.1111/biom.12027>
- Broms, K. M. (2013). Using Presence-Absence Data on Areal Units to Model the Ranges and Range Shifts of Select South African Bird Species. PhD thesis.
- Broms, K. M., Johnson, D. S., Altwegg, R., & Conquest, L. L. (2014). Spatial occupancy models applied to atlas data show Southern Ground Hornbills strongly depend on protected areas. *Ecological Applications*, 24(2), 363–374. <https://doi.org/10.1890/12-2151.1>
- Choi, H. M., & Hobert, J. P. (2013). The polygamma gibbs sampler for bayesian logistic regression is uniformly ergodic. *Electronic Journal of Statistics*, 7, 2054–2064. <https://doi.org/10.1214/13-EJS837>
- Clark, A. E., Altwegg, R., & Ormerod, J. T. (2016). A variational Bayes approach to the analysis of occupancy models. *PLoS ONE*, 11(2), e0148966. <https://doi.org/10.1371/journal.pone.0148966>
- Dorazio, R. M., & Rodriguez, D. T. (2012). A Gibbs sampler for Bayesian analysis of site-occupancy data. *Methods in Ecology and Evolution*, 3(6), 1093–1098. <https://doi.org/10.1111/j.2041-210X.2012.00237.x>
- Drouilly, M., Clark, A., & O’Riain, M. J. (2018). Multi-species occupancy modelling of mammal and ground bird communities in rangeland in the karoo: A case for dryland systems globally. *Biological Conservation*, 224, 16–25. <https://doi.org/10.1016/j.biocon.2018.05.013>
- Gardner, C. L., Lawler, J. P., Ver Hoef, J. M., Magoun, A. J., & Kellie, K. A. (2010). Coarse-scale distribution surveys and occurrence probability modeling for wolverine in interior Alaska.

- Journal of Wildlife Management*, 74(8), 1894–1903. <https://doi.org/10.2193/2009-386>
- Gelfand, A. E., Schmidt, A. M., Wu, S., Silander, J. A., Latimer, A., & Rebelo, A. G. (2005). Modelling species diversity through species level hierarchical modelling. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 54(1), 1–20.
- Gelfand, A. E., & Vounatsou, P. (2003). Proper multivariate conditional autoregressive models for spatial data analysis. *Biostatistics*, 4(1), 11–15. <https://doi.org/10.1093/biostatistics/4.1.11>
- Geweke, J. (1992). Evaluating the accuracy of sampling-based approaches to the calculations of posterior moments. *Bayesian Statistics*, 4, 641–649.
- Hanks, E. M., Schliep, E. M., Hooten, M. B., & Hoeting, J. A. (2015). Restricted spatial regression in practice: geostatistical models, confounding, and robustness under model misspecification. *Environmetrics*, 26(4), 243–254. <https://doi.org/10.1002/env.2331>
- Hodges, J. S., & Reich, B. J. (2010). Adding spatially-correlated errors can mess up the fixed effect you love. *The American Statistician*, 64(4), 325–334. <https://doi.org/10.1198/tast.2010.10052>
- Hoeting, J. A., Leecaster, M., & Bowden, D. (2000). An improved model for spatially correlated binary responses. *Journal of Agricultural, Biological, and Environmental Statistics*, 5(1), 102–114. <https://doi.org/10.2307/1400634>
- Hoffman, M. D., & Gelman, A. (2014). The No-U-turn sampler: Adaptively setting path lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research*, 15(1), 1593–1623.
- Holmes, C. C., & Held, L. (2006). Bayesian auxiliary variable models for binary and multinomial regression. *Bayesian Analysis*, 1(1), 145–168. <https://doi.org/10.1214/06-BA105>
- Hooten, M. B., Larsen, D. R., & Wikle, C. K. (2003). Predicting the spatial distribution of ground flora on large domains using a hierarchical bayesian model. *Landscape Ecology*, 18(5), 487–502. <https://doi.org/10.1023/A:1026001008598>
- Hughes, J., & Haran, M. (2013). Dimension reduction and alleviation of confounding for spatial generalized linear mixed models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 75(1), 139–159. <https://doi.org/10.1111/j.1467-9868.2012.01041.x>
- Huntley, B., Collingham, Y. C., Green, R. E., Hilton, G. M., Rahbek, C., & Willis, S. G. (2006). Potential impacts of climatic change upon geographical distributions of birds. *Ibis*, 148(s1), 8–28. <https://doi.org/10.1111/j.1474-919X.2006.00523.x>
- Hutchinson, R. A., Valente, J. J., Emerson, S. C., Betts, M. G., & Dietterich, T. G. (2015). Penalized likelihood methods improve parameter estimates in occupancy models. *Methods in Ecology and Evolution*, 6(8), 949–959. <https://doi.org/10.1111/2041-210X.12368>
- Johnson, D. S., Conn, P. B., Hooten, M. B., Ray, J. C., & Pond, B. A. (2013). Spatial occupancy models for large data sets. *Ecology*, 94(4), 801–808. <https://doi.org/10.1890/12-0564.1>
- Kellner, K. (2014). jagsui: Run JAGS (specifically, libjags) from R: an alternative user interface for rjags. *R package version*, 1.
- Kelsall, J., & Wakefield, J. (1999). Contribution to: "Bayesian models for spatially correlated disease and exposure data". In N. G. Best, L. A. Waller, A. Thomas, E. M. Conlon & R. Arnold (Eds.), *Bayesian Statistics 6, Proceedings of the Sixth Valencia International Meeting*, (Vol. 6 pp. 51). Oxford, UK: Oxford University Press.
- Latimer, A. M., Wu, S., Gelfand, A. E., & Silander, J. A. (2006). Building statistical models to analyze species distributions. *Ecological Applications*, 16(1), 33–50. <https://doi.org/10.1890/04-0609>
- Lichstein, J. W., Simons, T. R., Shriner, S. A., & Franzreb, K. E. (2002). Spatial autocorrelation and autoregressive models in ecology. *Ecological Monographs*, 72(3), 445–463. [https://doi.org/10.1890/0012-9615\(2002\)072\[0445:SAAAM\]2.0.CO;2](https://doi.org/10.1890/0012-9615(2002)072[0445:SAAAM]2.0.CO;2)
- Link, W. A., & Eaton, M. J. (2012). On thinning of chains in MCMC. *Methods in Ecology and Evolution*, 3(1), 112–115. <https://doi.org/10.1111/j.2041-210X.2011.00131.x>
- MacKenzie, D. I., Nichols, J. D., Lachman, G. B., Droege, S., Andrew Royle, J., & Langtimm, C. A. (2002). Estimating site occupancy rates when detection probabilities are less than one. *Ecology*, 83(8), 2248–2255. [https://doi.org/10.1890/0012-9658\(2002\)083\[2248:ESORWD\]2.0.CO;2](https://doi.org/10.1890/0012-9658(2002)083[2248:ESORWD]2.0.CO;2)
- Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., & Teller, E. (1953). Equation of state calculations by fast computing machines. *The Journal of Chemical Physics*, 21(6), 1087–1092. <https://doi.org/10.1063/1.1699114>
- Moreno, M., & Lele, S. R. (2010). Improved estimation of site occupancy using penalized likelihood. *Ecology*, 91(2), 341–346. <https://doi.org/10.1890/09-1073.1>
- Nelder, J., & Wedderburn, R. (1972). Generalized linear model. *Journal of the Royal Statistical Society*, 135(3), 370–384. <https://doi.org/10.2307/2344614>
- Plummer, M. (2003). JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling. In *Proceedings of the 3rd international workshop on distributed statistical computing*, (Vol. 124, pp. 125). Wien, Austria: Technische Universit at Wien.
- Plummer, M., Best, N., Cowles, K., & Vines, K. (2006). CODA: Convergence diagnosis and output analysis for MCMC. *R News*, 6(1), 7–11.
- Polson, N. G., Scott, J. G., & Windle, J. (2013). Bayesian inference for logistic models using Pólya-Gamma latent variables. *Journal of the American Statistical Association*, 108(504), 1339–1349. <https://doi.org/10.1080/01621459.2013.829001>
- Reich, B. J., Hodges, J. S., & Zadnik, V. (2006). Effects of residual smoothing on the posterior of the fixed effects in disease-mapping models. *Biometrics*, 62(4), 1197–1206. <https://doi.org/10.1111/j.1541-0420.2006.00617.x>
- Robert, C., & Casella, G. (2013). *Monte Carlo statistical methods*. New York: Springer-Verlag.
- Royle, J. A., & Dorazio, R. M. (2008). *Hierarchical modeling and inference in ecology: The analysis of data from populations, metapopulations and communities*. San Diego, CA: Academic Press.
- Tanner, M. A., & Wong, W. H. (1987). The calculation of posterior distributions by data augmentation. *Journal of the American Statistical Association*, 82(398), 528–540. <https://doi.org/10.1080/01621459.1987.10478458>
- Vieilledent, G., Merow, C., Guélat, J., Latimer, A., Kéry, M., Gelfand, A., ... Silander, J. Jr (2014). hsdm: Hierarchical bayesian species distribution models. R package version 1.4.
- Waller, L. A., & Gotway, C. A. (2004). *Applied spatial statistics for public health data*, Vol. 368. New York: Wiley.

SUPPORTING INFORMATION

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

How to cite this article: Clark AE, Altwegg R. Efficient Bayesian analysis of occupancy models with logit link functions. *Ecol Evol*. 2019;9:756–768. <https://doi.org/10.1002/ece3.4850>

APPENDIX

Certain posterior distributions

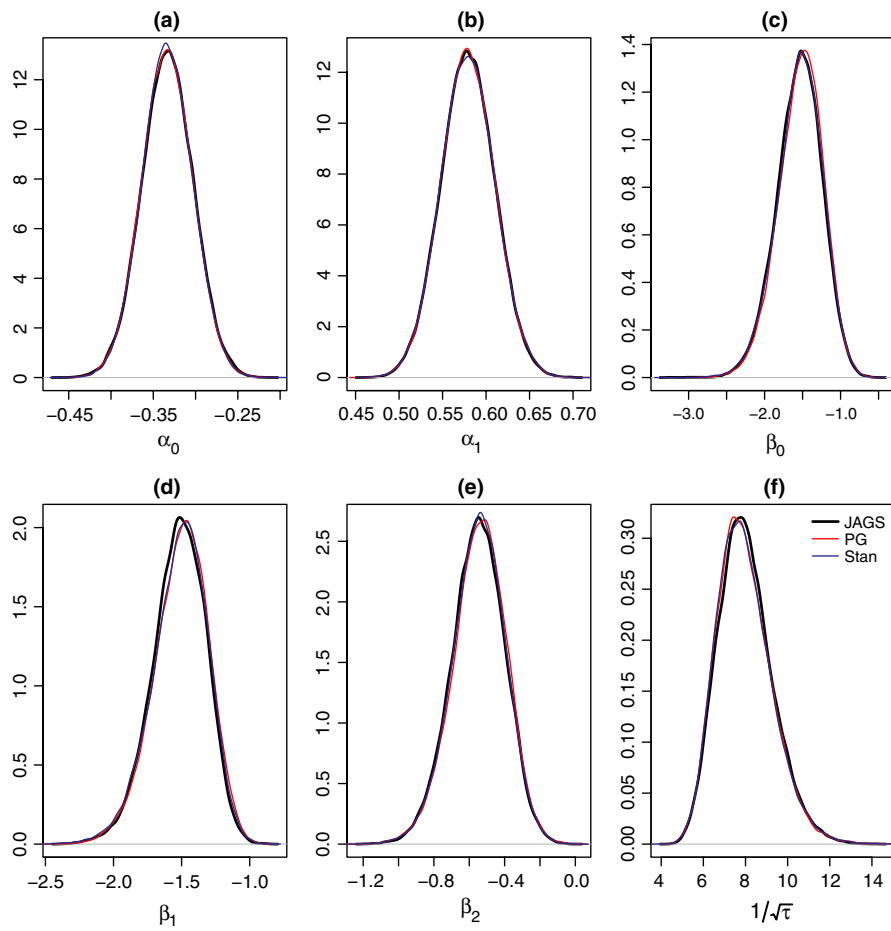


Figure A1. Posterior distributions of the parameters of the Bayesian spatial occupancy model using JAGS, Stan, and the Pólya-Gamma formulation for the Cape weaver data set [(a) = α_0 , (b) = α_1 , (c) = β_0 , (d) = β_1 , (e) = β_2 , (f) = $\frac{1}{\sqrt{r}}$]

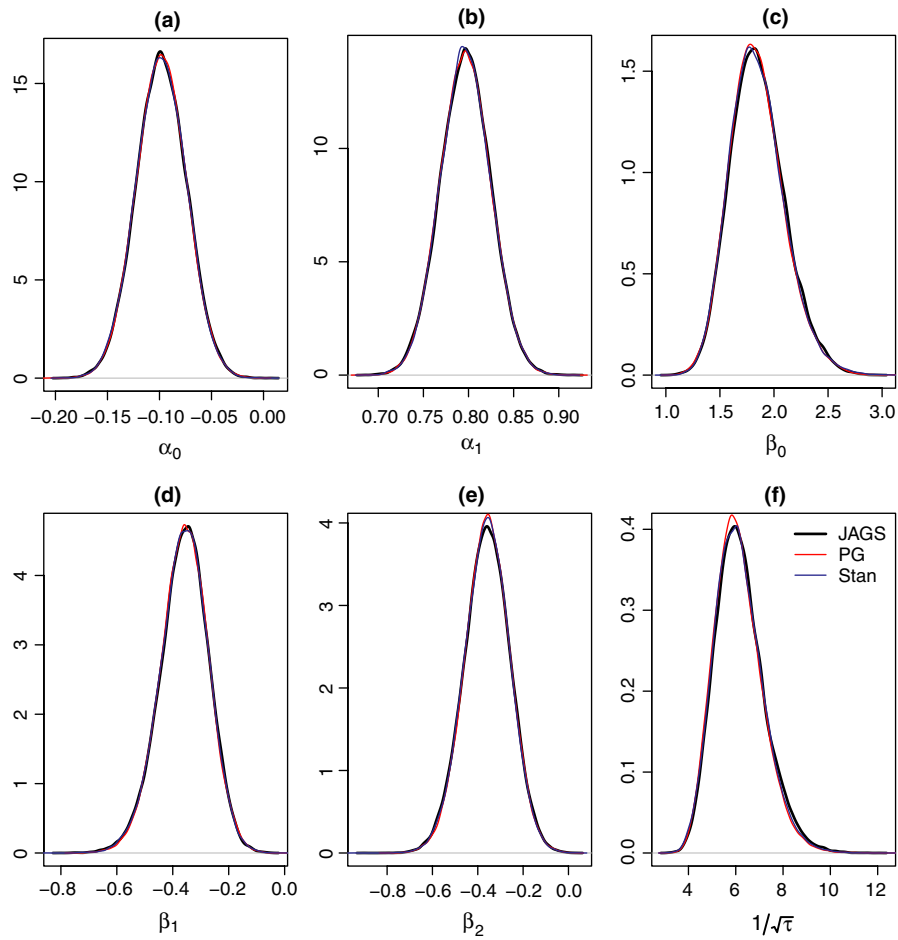


Figure A2. Posterior distributions of the parameters of the Bayesian spatial occupancy model using *JAGS*, *Stan*, and the Pólya-Gamma formulation for the helmeted guineafowl data set [(a) = α_0 , (b) = α_1 , (c) = β_0 , (d) = β_1 , (e) = β_2 , (f) = $\frac{1}{\sqrt{\tau}}$]

APPENDIX

Certain lagged sample autocorrelation functions

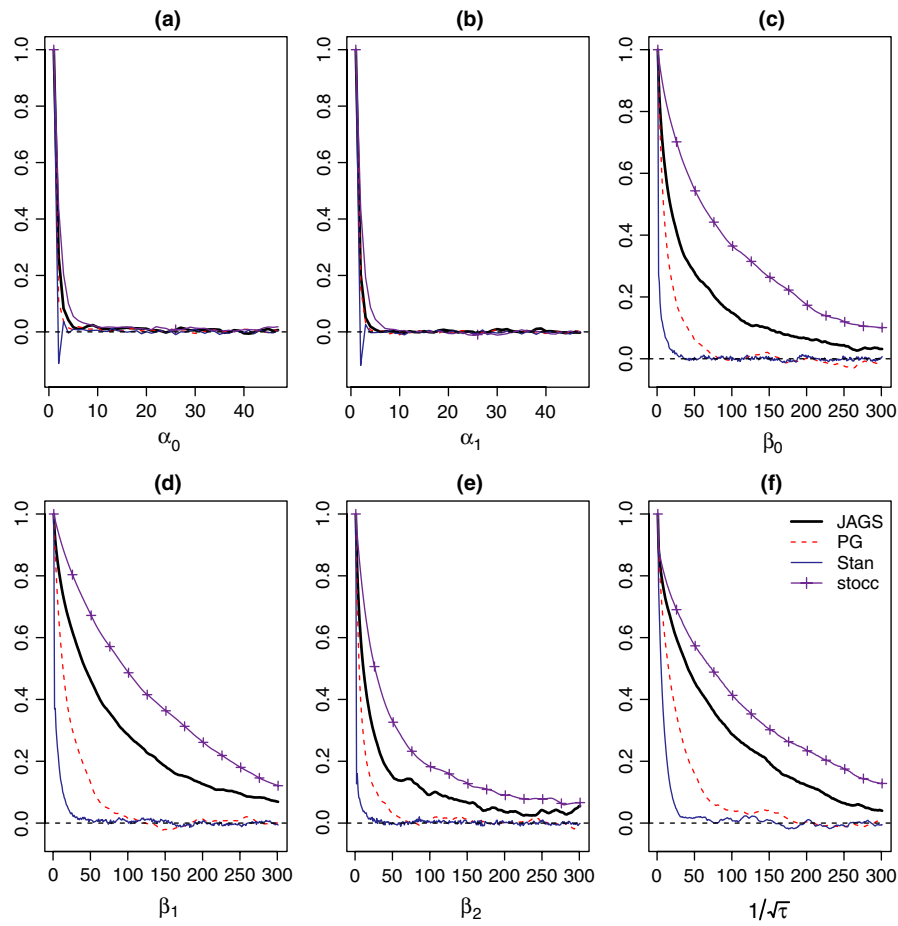


Figure A3. Estimated lagged sample autocorrelations of the posterior samples of the parameters of the Bayesian spatial occupancy model using *JAGS*, the Pólya-Gamma formulation, *Stan* and *stocc* for the Cape weaver data set [(a) = α_0 , (b) = α_1 , (c) = β_0 , (d) = β_1 , (e) = β_2 , (f) = $\frac{1}{\sqrt{r}}$]

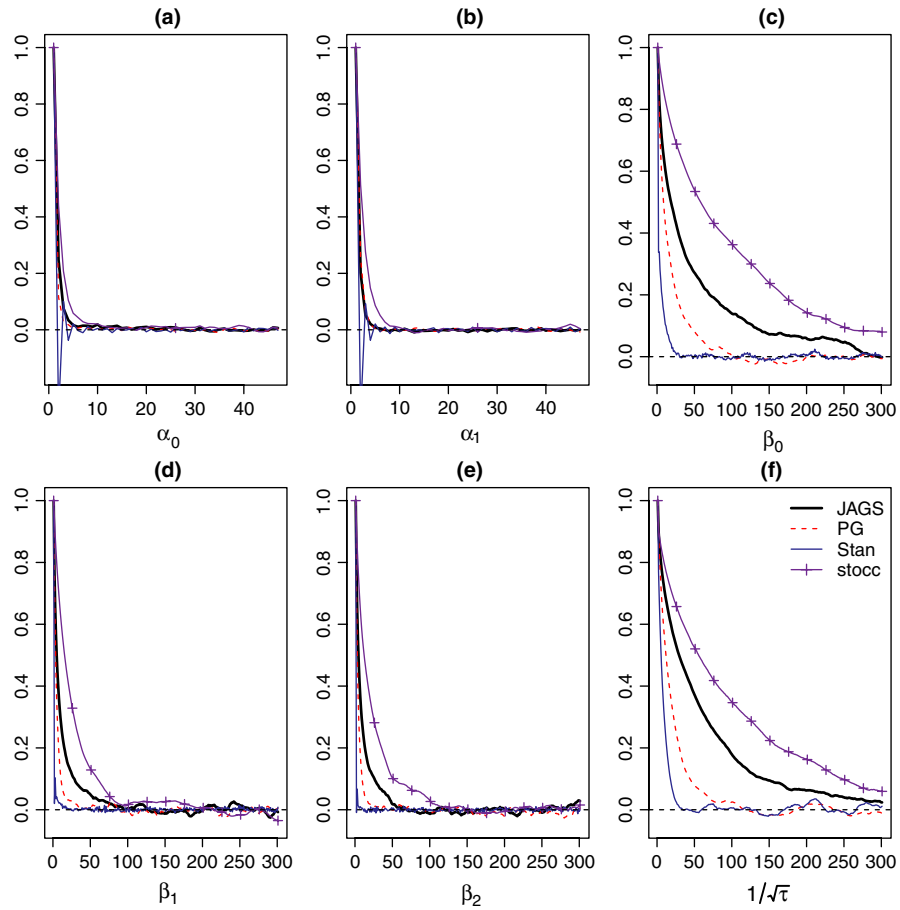


Figure A4. Estimated lagged sample autocorrelations of the posterior samples of the parameters of the Bayesian spatial occupancy model using JAGS, the Pólya-Gamma formulation, *Stan* and *stocc* for the helmeted guineafowl data set [(a) = α_0 , (b) = α_1 , (c) = β_0 , (d) = β_1 , (e) = β_2 , (f) = $\frac{1}{\sqrt{\tau}}$]