# Zero-inflated Beta distribution regression modeling

Running head: Zero-inflated Beta regression

Becky Tang[1*], Henry A. Frye[2], Alan E. Gelfand[1], John A Silander, Jr.[2]

[1]Department of Statistical Science, Duke University, Durham, NC 27708, USA; and [2]Department of Ecology and Evolutionary Biology, University of Connecticut, Storrs, CT 06269, USA

* Correspondence author: becky.tang@duke.edu

## Abstract

A frequent challenge encountered with ecological data is how to interpret, analyze, or model data having a high proportion of zeros. Much attention has been given to zero-inflated count data, whereas models for non-negative continuous data with an abundance of 0s are much fewer. We consider zero-inflated data on the unit interval and provide modeling to capture two types of 0s in the context of a Beta regression model. We model 0s due to missing by chance through left censoring of a latent regression, and 0s due to unsuitability using an independent Bernoulli specification. We extend the model by introducing spatial random effects. We specify models hierarchically, employing latent variables, and fit them within a Bayesian framework. Our motivating dataset consists of percent cover abundance of two plant families at a collection of sites in the Cape Floristic Region of South Africa. We find that environmental features enable learning about both types of 0s as well as positive percent cover. We also show that the spatial random effects model improves predictive performance. The proposed modeling enables ecologists to extract a better understanding of an organism's absence due to unsuitability vs. missingness by chance, as well as abundance behavior when present.

**Keywords**: Bayesian inference; Greater Cape Floristic Region; hierarchical model; hurdle model; percent cover; spatial random effects

# 1    Introduction

A frequent and continuing challenge encountered with ecological data is how to interpret, analyze, or model data having a high proportion of zeros. This *zero-inflation problem* (Martin et al. 2005) is found in scenarios such as: the binary presence/absence of individual organisms, counts of individuals, ranked abundances as in many vegetation plot datasets, and measures of biomass or proportion of area occupied by individuals in plot-level survey or remotely sampled data sets. A wealth of literature (e.g. Martin et al. 2005; Veech et al. 2016; Blasco-Moreno et al. 2019) asserts that zeros in a dataset can arise from multiple sources. Such prior studies have dichotomized zeros into "true" and "false" zeros; we avoid such value labels here, preferring stochastic interpretation.

While many statistical techniques have been proposed to address zero-inflation in discrete data, particularly for animal populations (Veech et al. 2016; Blasco-Moreno et al. 2019), less attention has been focused on continuous data. In particular, for count data, we remind the reader of the much-employed zero-inflated Poisson model (ZIP) (Lambert 1992) which introduces an additional point mass at zero. Models such as the zero-inflated negative binomial and other power series distributions, zero-inflated N-mixture models, and a general class of count transform models have also been developed (Hall 2000; Ghosh et al. 2006; Wenger and Freeman 2008; Veech et al. 2016; Siegfried and Hothorn 2020).

The zero-inflation challenge we focus on here is modeling the percent cover of plants within sampling plots yielding data on $[0, 1]$. Ecologists often assess plant community composition and abundance through visual assessments of percent cover, which has been shown to be a robust and repeatable measure of vegetation cover (Vanha-Majamaa et al. 2000). Percent cover is also often calculated by commonly used ordinal scales such as the Braun-Blanquet scale (van der Maarel 2007). In fact, our modeling can be applied to a variety of proportional cover data found in plant and animal surveys, e.g. marine invertebrates on hard substrate or coral reef communities (Bell and Galzin 1984), leaf cover estimates for plant phenological studies (Xie, Wang, Wilson and Silander 2018; Xie, Civco and Silander 2018), and biomass scaled as a proportion of area or organ mass (Jenkins et al. 2003). Absence, i.e., a zero, can be assumed to arise from one of three sources: unsuitability, random chance, or detection error. Unsuitability results from the biotic and abiotic factors impacting plant populations. If an area is suitable for a plant species but not found, we assume that this is due to random chance (i.e. it has not yet dispersed to that site). The explanation of failed detection of presence is not pursued here. In general, it requires detection/non-detection data from more than one independent visit to a sample unit. Further, it is extremely unlikely in our example system which is comprised almost exclusively of evergreen perennial species in an open shrubland. Lastly, we do not address the notion of "1-inflation" here.

Percent cover provides a challenge in that the data is continuous on $[0, 1]$. Though Beta

regression (Ferrari and Cribari-Neto 2004) is a customary choice to model such data (Douma and Weedon 2019), it has only been recently elaborated in detail for the context of vegetation percent cover (Damgaard and Irvine 2019). However, as we clarify below, these models are not adequate for what is truly zero-inflated, continuous data.

The key issue here is that discrete distributions explicitly place point mass at 0 while continuous distributions do not. So, a different mechanism is required to create mass at 0. One method of introducing zeros for continuous data on $[0, 1]$ is through the left-censoring of a latent random variable at 0. For example, consider a continuous variable $W$ with mass on $(-1, 1)$. Assume that the actual observation is a left-censored at 0 realization of this variable. Then, the probability of observing a 0 becomes $P(W \leq 0)$. This approach is similar to the Tobit model which has been long used in economics (Amemiya 1984) and applied in remote sensing (Peterson 2005). The Tobit model introduces a regression for latent Gaussian variables and creates point mass at 0 through left-censoring at 0 (Chib 1992; Long and Long 1997).

A different approach to modeling zero-inflated data is a hurdle model (Mullahy 1986). Hurdle models explain a chosen measure of abundance given presence. In our setting, they model the positive percent coverages, ignoring the observed 0's. Separately, the observation at each site is given a binary response according to presence or absence, modeled through a binary regression. This zero-inflation for a beta distribution, using site-specific independent Bernoulli variables, has been termed a *zero-inflated Beta* model in Ospina and Ferrari (2012) and Ospina and Ferrari (2010). They refer to this hurdle model as BEZI (BEta Zero Inflated). For observations on $[0, 1)$, neither of the foregoing models (BEZI or left-censored) is really a zero-inflated continuous data specification. In both cases, a single mass at 0 is introduced. We build upon the previous literature and propose a new model that accommodates two sources of 0 for percent cover, analogous to the zero-inflated Poisson model for counts. .

For interpretation, the left-censored regression type of 0 will be referred to as absence by chance (Blasco-Moreno et al. 2019). The Bernoulli point mass at 0 will be referred to as absence due to unsuitability. Our contribution here is to formalize a zero-inflated beta regression model for percent cover. This entails a regression within a left censored Beta specification to capture absence by chance as well as a regression to capture absence due to unsuitability. Can we separate both sources of 0's? As we discuss below, the answer is yes given the use of informative priors for identifiability.

As the environmental covariates may not be sufficient to explain all the variability in the data, introducing spatially correlated random effects is expected to provide improved model performance. Therefore, employing the geo-coded locations of the sampling sites, we offer a spatial random effects version of the foregoing specification. Spatial zero-inflated models have been proposed in the literature (e.g. Agarwal et al. 2002; Rathbun and Fei 2006) and we extend this work to our proposed zero-inflated Beta regression model. This raises the

question of model comparison; is the spatial model preferred? We address this through novel model comparison for these zero-inflated distributions. As a case study, we consider two plant families (the Restionaceae and Crassulaceae) located across a suite of sample plots in the Cape Floristic Region (CFR) in South Africa.

The format of the paper is as follows: Section 2 describes the CFR dataset which motivates our modeling. Section 3 presents modeling details for the nonspatial case, offers model fitting details as well as model assessment and comparison approaches, and finishes with the spatial model. Section 4 offers simulation investigations to serve as a proof of concept and compare various models. Section 5 presents the results of analyses for the CFR data.
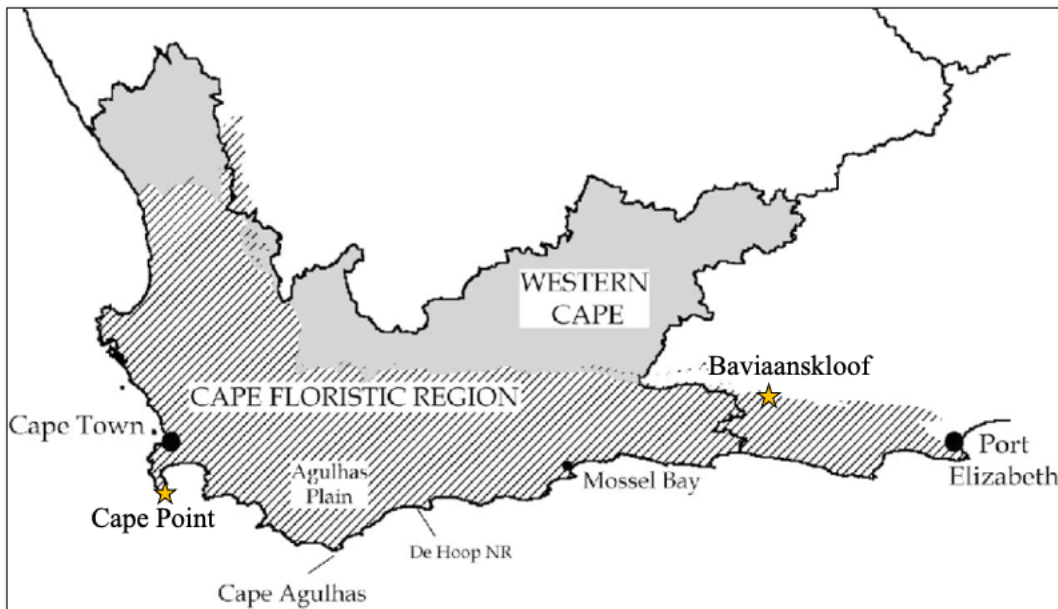
# 2 The dataset

We consider percent cover for species within two plant families: Crassulaceae and Restionaceae in the Cape Floristic Region (CFR). The Crassulaceae found here are typically perennial succulents that contribute to the CFR's diversity of succulent flora (Born et al. 2006). They are arid adapted plants; we would expect them to inhabit drier and warmer microhabitats within the CFR, likely in sites protected by fire. The Restionaceae are evergreen reed-like plants that are an iconic component of the CFR. They are found in a wide range of habitats across the nutrient poor soils of the CFR (Linder 2001).

Our data were collected between 2010 and 2011 from plots concentrated in two regions of South Africa: 61 plots in the Cape Point section of Table Mountain National Park and 119 plots in the Baviaanskloof Mega-Reserve (Fig. 1). Following the relevé method (Ellenberg and Mueller-Dombois 1974), a well-established inventory protocol used to classify vegetation and inventory species (e.g Schaminée et al. 2009; Dengler et al. 2011), each 10 x 5 m plot was sampled for its floristic composition using observations of percent cover. Histograms of the observed percents are presented in Fig. 2a. 49% of the observations for Crassulaceae and 31% of observations for Restionaceae are 0, arguing for our proposed introduction of zero-inflation.

The plots are associated with a variety of environmental variables georeferenced for each plot and largely derived from the WorldClim (Hijmans et al. 2005) and Schulze (Schulze 1997) climate datasets. We employ four of these variables in our regression modeling: mean annual evapotranspiration (1-minute grid cell), mean annual precipitation (30 arc second), average minimum temperature for July (the austral winter; 1-minute grid cell), rainfall concentration index (1-minute grid cell; scaled 0 (rainfall is evenly divided across all months) to 100 (rain only falls one month of the year)). Fig. S1 displays histograms of the environmental covariates, and Fig. 2b plots the positive percent covers of both species against the environmental covariates as well as the distribution of environmental covariates for observed

4

Figure 1: Location of study regions Cape Point and Baviaanskloof within the Cafe Floristic Region (map adapted from Turpie et al. (2003)).



0s. As our observations are spatially referenced, we may also be interested in investigating if the data are spatially correlated for a given family. The percent covers in Baviaanskloof and Cape Point plotted by longitude-latitude are presented in Figs 2c and S2, respectively. Throughout this work, 'S' denotes Supplementary Information.

# 3 Building the model

As we have discussed above, we propose a site-level model that will have a Beta distribution specification for the continuous observations. It will employ left censoring to create a point mass at 0 which corresponds to missingness by chance at the site. It will introduce a Bernoulli trial at the site to capture the probability of unsuitability.

## 3.1 Left-censored Beta regression

We first introduce a 0 through left censoring. In a general setting, for site $i$, let the observation, $Y_i = \max(0, W_i)$ where $W_i$ is a continuous variable on $(-\infty, \infty)$ with distribution $f_i$ so $P(Y_i = 0) = P(W_i \leq 0) = \int_{-\infty}^{0} f_i(w)dw$. When $f_i$ is a normal distribution this model is a

Tobit model (Tobin 1958). In our context, we take $f_i$ to be the Beta distribution. However, the Beta has bounded support on $(0, 1)$. Therefore, we propose extending the Beta distribution: if $V \sim \text{Beta}(\alpha, \beta)$, then for $0 < a < \infty$, $W = (a + 1)V - a$ provides an extended or re-scaled Beta on $(-a, 1)$. The probability of a negative $W$ provides our Beta mass at 0: $P(W \leq 0) = P(V \leq \frac{a}{a+1} \equiv c)$ (Fig. 3).

We follow a Bayesian approach for modeling fitting. We introduce a latent Beta variable $V_i$ such that $Y_i = \max\{0, W_i = (a + 1)V_i - a\}$. Then for all $Y_i > 0$, we immediately know $V_i$. For $Y_i = 0$, the associated $V_i$ is less than or equal to $c$ by construction. Thus, introducing the latent $V_i$ allows us to have the 'complete' data. In fitting the model, we work with a convenient reparametrization for Beta regression (Ferrari and Cribari-Neto 2004). We employ $\mu_i = \alpha_i/(\alpha_i + \beta_i)$, the mean parameter and $\nu_i = \alpha_i + \beta_i$, the so-called "sample size" parameter. The distribution of the extended Beta on $(-a, 1)$ is also parameterized in terms of $\mu_i$ and $\nu_i$. We model $logit(\mu_i) = \mathbf{X}_i^T \boldsymbol{\delta}$, and $log(\nu_i) = \mathbf{X}_i^T \boldsymbol{\psi}$ where $\mathbf{X}_i$ is a vector of environmental covariates associated with site $i$.

With regard to addressing the effect of choice of $a$, the Appendix offers some analytical investigation. The choice of $a = 1$ mirrors the Tobit model by providing matching support above and below 0. Ideally, we would take $a$ to be a model parameter and learn about its value during model fitting, along with $\boldsymbol{\delta}$ and $\boldsymbol{\nu}$. In implementation, this leads to identifiability issues: $a$ (or equivalently, $c$) and the intercept $\delta_0$ for $\mu$ compete to explain the probability of 0. One approach to address this issue would be to use an informative prior on $c$ or $\delta_0$. However, we remark that the value of $c$ does not have interpretation; it is a mechanism used to obtain 0's. Practically, this suggests holding $c$, equivalently $a$, fixed at a pre-specified value during model fitting and then doing out-of-sample model comparison.

In this regard, and following the Appendix, through simulations we briefly explore sensitivity of inference to the choice of $a$. The simulations are performed by generating data under a $c_{\text{true}}$, fitting the censored model with several different values of $c_{fixed}$, and examining how recovery of the covariate effects and probabilities of 0 are affected. The regression for $\mu$ is $logit(\mu) = \delta_0 - 0.5x_1$, with different $\delta_0$ and $\psi$ in each simulation. We evaluate credible interval coverage of $\delta_1$ and root mean square predictive error (RMSPE) of the probabilities of 0 across 100 simulations for each pair $(c_{\text{true}}, c_{\text{fixed}})$. We find that generally, recovery of the probability of 0 is robust to the choice of fixed $a$; predicted probabilities on the held-out test set tend to yield similar RMSPEs (Fig. S3b). The largest differences in predictions and RMSPEs occur when $c_{\text{true}}$ is small and $c_{\text{fixed}}$ is large, or vice versa (Fig. S5).

Empirical coverage of the credible intervals for $\delta_1$ achieves or is slightly below nominal when $c$ is fixed at the true value, but tends to decrease as the difference between $c_{\text{true}}$ and $c_{\text{fixed}}$ increases (Fig. 3a). This decreasing coverage is exacerbated when data are generated with smaller intercept $\delta_0$, holding everything else fixed (ex. Simulations 1-3, Simulations 4-6). This is due to the larger number of 0s in the data induced by the smaller intercept.

However, in almost all pairs of $(c_{\text{true}}, c_{\text{fixed}})$, the sign of $\delta_1$ is correctly estimated (Fig. S4). Therefore, even if $c$ is fixed at an incorrect value, the model is expected to recover the correct covariate relationships.

From the simulations, unsurprisingly, we find that fixing $c$ at the truth leads to the best recovery of parameters and predicted probabilities of 0. We generally find that fixing $c$ to be large when the data contain a high proportion of 0s leads to better performance, and similarly for small $c$ with a smaller proportion of 0s. Inference is the least robust when the magnitude of the difference between $c_{\text{true}}$ and $c_{\text{fixed}}$ is large. Thus, if the modeler seeks to choose a single value *a priori*, perhaps fixing $c = 0.5$ ($a = 1$) is the most natural option. Again, we suggest analyzing the data with differences choices of $c$ and performing sensitivity analysis to determine the best fitting model. For the remainder of the text, we refer to this left-censored extended Beta with single-zero model as 'LCEB'.

## 3.2   Zero-inflated Beta modeling - the nonspatial case

Let $Z_i$ be a Bernoulli variable with $P(Z_i = 1) = \pi_i$ and let $f_i$ be a density on $(0, 1)$. Then, suppose $Y_i = 0$ if $Z_i = 1$, and $Y_i \sim f_i(y)$ if $Z_i = 0$. Here, $P(Y_i = 0) = \pi_i$. This is the BEZI model of Ospina and Ferrari (2010). It is a hurdle model analogous to the setting where $Y_i$ is a count variable and $f_i$ is the Poisson distribution truncated at 1. (Mullahy 1986).

If $Z_i = 0$, then let $Y_i$ be a realization of the censored Beta regression model described in Section 3.1: $Y_i = \max(0, W_i)$ where $W_i = (a + 1)V_i - a$ and, using the Ferrari and Cribari-Neto (2004) parameterization, $V_i \sim Beta(\mu_i, \nu_i)$. Now, $P(Y_i = 0) = P(Z_i = 1) + P(Z_i = 0 \cap W_i \leq 0) = \pi_i + (1 - \pi_i) \int_{-a}^{0} f_i(w)dw$ (so $P(Y_i > 0) = (1 - \pi) \int_0^1 f_i(w)dw$). We immediately see the two sources of 0's as well as the parallel construction with the familiar zero-inflated count model. In this regard, we will refer to $(1 - \pi_i) \int_{-1}^{0} f_i(w)dw$ as the probability of absence by chance with $\pi_i$ the inflated probability of absence which we ascribe to unsuitability.

We can identify the two sources of 0's if $a$ is held fixed because the $\mu_i$ and $\nu_i$ are informed by all of the data while only the 0's inform about $\pi_i$. However, informative prior information is needed to capture the relative magnitudes of the chances of each source. In the absence of further knowledge, this is addressed via out-of-sample model comparison. We refer to our proposed zero-inflated Beta model as $\mathcal{M}_1$.

Again, following the above, for all $Y_i > 0$, we immediately know $V_i$. For $Y_i = 0$, we know that if the observation arose from the Beta, by construction, the associated $V_i$ is less than or equal to $\frac{a}{1+a} = c$. Our zero-inflated model takes the form

$$\begin{aligned}
P(Y_i = 0 | \pi_i, \mu_i, \nu_i) &= \pi_i + (1 - \pi_i)P(V_i \leq c | \mu_i, \nu_i) \\
f(y_i | \pi_i, \mu_i, \nu_i) &= (1 - \pi_i)f_{v_i}(v_i | \mu_i, \nu_i), \quad y_i > 0
\end{aligned} \tag{1}$$

Here, a latent indicator variable $Z_i$ is introduced to determine the source of an observation: if $Z_i = 1$ then the associated $Y_i$ arises as a 0 from the degenerate process, and if $Z_i = 0$ then $Y_i$ is a realization of the extended Beta. Then we have that $P(Y_i = 0, Z_i = 1|\pi_i, \mu_i, \nu_i) = \pi_i$ and $P(Y_i = 0, Z_i = 0|\pi_i, \mu_i, \nu_i) = (1 - \pi_i)P(V_i \leq c|\mu_i, \nu_i)$.

Given the parameters $\pi_i, \mu_i, \nu_i$ and the latent $Z_i$, the likelihood is:

$$L(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\nu}, \mathbf{Z}; \mathbf{Y}) \quad = \prod_{i=1}^{n} \pi_i^{z_i} \cdot ((1 - \pi_i)f_{y_i}(y_i|\mu_i, \nu_i))^{1-z_i}. \tag{2}$$

We model logit$\pi_i = \mathbf{G}_i^T \boldsymbol{\gamma}$, and as in the case of the single-zero censored Beta, logit$\mu_i = \mathbf{X}_i^T \boldsymbol{\delta}$, and $\log \nu_i = \mathbf{X}_i^T \boldsymbol{\psi}$ where $\mathbf{X}_i, \mathbf{G}_i$. The context will determine the commonalities between $\mathbf{X}_i$ and $\mathbf{G}_i$ with regard to explaining $\pi$ and the parameters of the Beta distribution. Then, we wish to infer about the regression coefficients, $\boldsymbol{\gamma}, \boldsymbol{\delta}, \boldsymbol{\nu}$, as well as the latent $Z_i$ associated with $Y_i = 0$. With priors, the posterior distribution is proportional to

$$\prod_{i=1}^{n} f_{Y_i|Z_i, \boldsymbol{\delta}, \boldsymbol{\nu}}(Y_i|Z_i, \boldsymbol{\delta}, \boldsymbol{\nu}) \cdot \prod_{i=1}^{n} f_{Z_i|\boldsymbol{\gamma}}(Z_i|\boldsymbol{\gamma}) \cdot f_{\boldsymbol{\gamma}}(\boldsymbol{\gamma}) \cdot f_{\boldsymbol{\delta}}(\boldsymbol{\delta}) \cdot f_{\boldsymbol{\psi}}(\boldsymbol{\psi}). \tag{3}$$

Prior choices follow in the next subsection.

## 3.3   Model fitting

We use a Gibbs sampler with Metropolis updates for the $\boldsymbol{\delta}, \boldsymbol{\gamma}, \boldsymbol{\psi}$. If $Y_i > 0$, we immediately have $Z_i$ and $V_i$. Then for all $i$ where $Y_i = 0$, at every iteration within the sampler we sample $Z_i$ from its Bernoulli full conditional with success probability $P(Z_i = 1|Y_i = 0) = \pi_i/(\pi_i + (1 - \pi_i)P(V_i \leq c))$. If $Z_i = 0$, we sample an associated latent $V_i$ from its conditional Beta distribution restricted to the interval $(0, c)$, analogous to sampling from truncated Normal variables in Bayesian Tobit modeling (Chib 1992). Given $\mathbf{Z}$ and $\mathbf{V}$, the parameters $\boldsymbol{\gamma}$ are independent of the parameters $(\boldsymbol{\delta}, \boldsymbol{\psi})$. In the sequel, we assume $\nu_i = \nu$, i.e., a common sample size parameter across sites. This is in accord with Damgaard and Irvine (2019) who use the parameterization, $\mu_i$ with common shape parameter $\delta = 1/(\nu + 1)$.

Distinguishing the two types of zeros in the data is challenging. The magnitudes of $\pi_i$ and $\mu_i$ are strongly influenced by their intercepts $\gamma_0$ and $\delta_0$, respectively. A large positive $\gamma_0$ will result in $\pi_i$ close to one, encouraging unsuitability absence; a large negative $\delta_0$ will result in $\mu_i$ small, encouraging absence by chance. Recall that the probability of a zero at location $i$ is $\pi_i + (1 - \pi_i)P(V_i \leq c|\mu_i, \nu_i)$. The second term leads to difficulty in separating the source of a zero due to weak identifiablity of the intercepts. Rather than removing an intercept from $\pi_i$ or $\mu_i$, we choose to adopt a strong prior on one of the intercepts.

This leads to two questions: which intercept, and what prior? The natural choice is

the intercept $\gamma_0$ for the incidence of inflated zeros; an ecologist may have a general belief of how (un)favorable a location is for a given species, but may have less knowledge about the expected percent cover. If so, a tight prior centered at the modeler's belief for $\gamma_0$ can be used. We take this approach in our simulations. However, we suggest performing modeling comparison or sensitivity analysis with different priors to select a model. Our proposed metrics for modeling comparison follow in the next subsection. We use non-informative or weakly-informative priors for all the remaining parameters. For a full list of priors used in all models, please see Table S1 in the Supplement.

## 3.4   Model assessment and comparison approaches

If the source of a zero is known–as with the simulated data below–one measure of model assessment is through classification of the type of zero. We use receiver operating characteristic (ROC) curves and the corresponding area under the curve (AUC) for the true source of a zero versus the predicted probability of arising from that source. In practice, we won't know the true source of a zero. Instead, we assess and compare models in terms of their ability to distinguish zero and positive observations. Another metric is Tjur's $R^2$ statistic, a coefficient of discrimination often used to evaluate logistic regression prediction (Tjur 2009). This statistic is calculated as $R^2 = \bar{\hat{p}}_1 - \bar{\hat{p}}_0$, where $\bar{\hat{p}}_1$ and $\bar{\hat{p}}_0$ are the average of fitted values for successes and failures. Larger AUC and $R^2$ indicate superior classification performance.

The continuous ranked probability score (CRPS) which compares the entire predictive cumulative distribution function, $F(x)$, to an observed held-out data point (Gneiting and Raftery 2007) offers further model comparison. For a continuous CDF, $CRPS = \int (F(x) - \mathbf{1}(y < x))^2 dx$, where $y$ is the observed value. For a discrete distribution, the ranked probability score replaces the integral with a sum. A smaller (C)RPS indicates that the predictive distribution is more concentrated around the observation and is thus preferred. In the case of our proposed distribution, we have a continuous predictive CDF with a point mass at zero. Let $p_0$ denote the mass at 0. Then our score takes the form:

$$(C)RPS = \begin{cases} p_0^2 + \int_0^\infty (F(x) - \mathbf{1}(y < x))^2 dx, & y > 0 \\ (1 - p_0)^2 + \int_{-\infty}^0 (1 - F(x))^2 dx, & y = 0 \end{cases}$$

Thus, we can obtain the (C)RPS for the full distribution, denoted as $CRPS_f$. How well the predictive distribution approximates the positive observations may also be of interest. In this case, we omit the zero observations from the calculation and only use the positive data. We refer to this resulting score as $CRPS_h$. We use Monte Carlo integration to evaluate the integrals (Supplement S1).

Because evaluation using our ranked probability score requires separating observations

9

based on whether they are zero or positive, the concentration of the predictive distribution depends on the number of sources of zeros that are being modeled. For example, it would be inappropriate to use our $CRPS_f$ to compare our proposed model with the BEZI model. However, if we wish to perform sensitivity analysis to assess the effect of prior choice in our proposed model, then the use of $CRPS_f$ would be appropriate to compare models.

## 3.5 Zero-inflated Beta modeling - the spatial case

With observations recorded at geo-coded (point-referenced) locations, spatial dependence can be introduced using random effects. Within our zero-inflated Beta regression model, we can insert these spatial random effects in two places. The first is in the mean specification, e.g., $\text{logit}\mu(\mathbf{s}_i) = \mathbf{X}_i^T\boldsymbol{\delta} + \eta(\mathbf{s}_i)$, where $\eta(\mathbf{s}_i)$ is the random effect at the geo-coded location $\mathbf{s}_i$. In this case, the random effect impacts both the positive percent covers as well as the classification of the zero source. The second introduces a spatial random effect $\tau(\mathbf{s}_i)$ into the probability that a site is unfavorable, e.g., $\text{logit}\pi(\mathbf{s}_i) = \mathbf{G}_i^T\boldsymbol{\gamma} + \tau(\mathbf{s}_i)$. Then, space impacts the classification of the zero source. Given an unobserved set of locations and associated predictors, we can perform spatial prediction for percent covers at these locations.

We model the set of spatial random effects $\eta(\mathbf{s})$ using a Gaussian process centered at 0 with exponential covariance function, $C(\mathbf{s}_i, \mathbf{s}_i + h) = \sigma^2 e^{-\phi_\eta \|h\|}$. As shown by Zhang (2004), the product of the spatial variance $\sigma^2$ and decay term $\phi$ is identifiable, but the individual terms are not. With greater interest in learning about spatial variability , we set $\phi$ to be fixed during model fitting and estimate $\sigma^2$ (see Banerjee et al. (2014) in this regard). The posteriors for both spatial models are provided in Supplement S3.

We fit the $\eta(\mathbf{s})$ ($\tau(\mathbf{s})$) with an elliptical slice sampler (Murray et al. 2010). As mentioned previously, we hold the spatial decay parameter fixed and estimate the spatial variance $\sigma^2$ for identifiability (Zhang 2004). We take a weakly informative prior Inverse Gamma(0.5, 0.5) prior on $\sigma_\eta^2$ ($\sigma_\tau^2$). Introducing a random effect in the regression to create $\pi(\mathbf{s}_i)$ or $\mu(\mathbf{s}_i)$ may lead to instability when estimating the corresponding intercept. Therefore, we suggest using an informative prior for the intercept in the regression where the random effect is introduced, and performing model assessment and comparison to determine the best model. For the remainder of the paper, let $\mathcal{M}_2$ denote the model where the spatial random effect $\eta(\mathbf{s}_i)$ is inserted into the mean for the extended Beta, and $\mathcal{M}_3$ the model where the random effect $\tau(\mathbf{s}_i)$ is included in the probability of unsuitability. Including the spatial process in models $\mathcal{M}_2$ and $\mathcal{M}_3$ leads to increased computational time compared to the non-spatial $\mathcal{M}_1$, especially as the size of the data set increases (ex. about 0.5 vs. 3 minutes for $n = 300$ using a server node with one core Intel(R) Xeon(R) CPU E5-2680 v3 @ 2.50GHz; Fig. S7).

# 4  Simulation

## 4.1  $\mathcal{M}_1$: a proof of concept

We generate data under our proposed zero-inflated Beta model $\mathcal{M}_1$, and examine how well we recover the truth. Data are generated and models fit with $a = 1$. We include a $N(0, 1)$ covariate for $\pi_i$, and a separate $N(0, 1)$ variable for $\mu_i$ and have a constant $\nu_i = \nu$. For the intercept in unsuitability, we use an informative Normal prior centered at the true $\gamma_0$ with 0.25 standard deviation. We run the samplers for 7500 iterations after 2500 burn-in, and thin to retain every fifth sample. We examine traceplots and autocorrelation to assess convergence. All analyses were conducted using R (Version 3.6.1) (R Core Team 2013).

As noted above, the intercepts in $\pi_i$ and $\mu_i$ largely control the amount of zeros in the data, as well as the proportions of the sources of zeros: Simulation 1.1 has equal proportions, Simulation 1.2 has more zeros arising from the Beta, and Simulation 1.3 has more zeros arising due to unsuitability. Plots of ROC curves demonstrate that our model is able to perform much better than random guessing in separating the two types of zeros (Fig. S6).

## 4.2  Model comparison: nonspatial models

We compare BEZI, LCEB, and $\mathcal{M}_1$ by generating train and test data under each of these models, fitting all three models on the sets of train data, and predicting on the sets of test data. In the case of LCEB, and $\mathcal{M}_1$, $a$ is fixed at 1. We use Tjur's $R^2$ and AUC to evaluate the ability of each model to distinguish zero and positive observations, and $CRPS_h$ to compare each model's predictive distributions for the positive observations. Comparisons are performed for varying amounts of 0s in the data and, for data generated under $\mathcal{M}_1$, with different proportions for the sources of 0s by changing $\delta_0$ and $\gamma_0$. We perform 50 replications for each data-generating model and set of $\delta_0$ and $\gamma_0$, and present results averaged across the 50 runs. More details on the simulations are provided in the Supplement S2.

When the data-generating model is $\mathcal{M}_1$, the true model outperforms the other models using all three prediction metrics (Table 1). The higher $R^2$ and AUC indicate its superior ability in separating the positive observations from the zero observations. As evidenced by the lower $CRPS_h$, $\mathcal{M}_1$'s predictive distribution for positive observations is more concentrated around the truth in all scenarios. These findings reveal that the BEZI and LCEB models cannot capture all the 0's well enough when data are generated with two sources. That is, the regression for $\pi$ in BEZI cannot explain all the 0's, and the regressions for $\mu$ and $\nu$ in LCEB cannot accommodate both the positive percent covers and the large incidence of 0's.

When BEZI is the data-generating model, we see that model $\mathcal{M}_1$ performs very similar to and sometimes better than the true model (Table S3). In particular, when $\delta_0$ is large and

11

positive (i.e. the mean of the positive percent covers is large), the model $\mathcal{M}_1$ sometimes outperforms the true model. However, when $\delta_0$ is negative, the BEZI model offers prediction performance. This can be attributed to the dual role that $\delta_0$ plays in $\mathcal{M}_1$, as $\delta_0$ helps determine the mean of the positive percent covers and also contributes to the probability of 0. Under the BEZI model, $\delta_0$ only has a single role. Thus when $\delta_0$ is negative, the probability of 0 must increase under $\mathcal{M}_1$ if the model wishes to capture the lower positive percent covers.

If the left-censoring model LCEB generates the data, we find that $\mathcal{M}_1$ outperforms or performs just as well as the true model in almost all scenarios (Table S4). Additionally, the BEZI model often performs just as well. This is explained by the extra $\pi$ parameter in both the $\mathcal{M}_1$ and BEZI models, as it provides flexibility to accommodate the behavior of 0s generated under the censor-only model. That is, the likelihood under LCEB is $L(a, \boldsymbol{\delta}, \boldsymbol{\psi}; \{Y_i\})$, while under BEZI it is $L(\boldsymbol{\gamma}, \boldsymbol{\delta}, \boldsymbol{\psi}; \{Y_i\})$ and under $\mathcal{M}_1$ it is $L(a, \boldsymbol{\gamma}, \boldsymbol{\delta}, \boldsymbol{\psi}; \{Y_i\})$. For this reason, LCEB has poor performance for data generated under BEZI and $\mathcal{M}_1$, precisely because it has fewer parameters and thus less flexibility (Tables 1, S3).

Table 1: Simulation for data generated under the proposed zero-inflated Beta model $\mathcal{M}_1$. $R^2$ and AUC for classifying between zero and positive, and $CRPS_h$ for predictive distribution of positive observations for various simulations. Simulations vary by amount of zeros in the data, as well as the proportions of zeros arising from each source (% unsuitable, % random chance). Bold indicates best performer.

| | $R^2$ | | | AUC | | | $CRPS_h$ | | |
|---|---|---|---|---|---|---|---|---|---|
| (% deg. 0, % Beta 0) | BEZI | LCEB | $\mathcal{M}_1$ | BEZI | LCEB | $\mathcal{M}_1$ | BEZI | LCEB | $\mathcal{M}_1$ |
| 18%, 5% | 0.046 | 0.033 | **0.064** | 0.622 | 0.590 | **0.656** | 0.085 | 0.115 | **0.058** |
| 12%, 12% | 0.016 | 0.092 | **0.096** | 0.574 | 0.647 | **0.667** | 0.138 | 0.083 | **0.073** |
| 5%, 19% | 0.002 | 0.160 | **0.168** | 0.561 | 0.746 | **0.750** | 0.144 | 0.080 | **0.077** |
| 38%, 15% | 0.074 | 0.034 | **0.097** | 0.634 | 0.604 | **0.672** | 0.125 | 0.150 | **0.083** |
| 22%, 26% | 0.007 | 0.152 | **0.158** | 0.563 | 0.716 | **0.719** | 0.148 | 0.097 | **0.084** |
| 13%, 37% | 0.003 | 0.193 | **0.198** | 0.560 | 0.752 | **0.753** | 0.145 | 0.096 | **0.087** |
| 71%, 9% | 0.046 | 0.025 | **0.051** | 0.600 | 0.594 | **0.647** | 0.151 | 0.191 | **0.086** |
| 39%, 40% | 0.014 | 0.125 | **0.135** | 0.579 | 0.729 | **0.748** | 0.140 | 0.124 | **0.099** |
| 15%, 56% | 0.002 | 0.183 | **0.187** | 0.550 | 0.772 | **0.771** | 0.135 | 0.104 | **0.094** |

## 4.3 Model comparison: spatial vs nonspatial zero-inflated models

Again, $\mathcal{M}_1$ is the non-spatial zero-inflated model we developed in Section 3.2, and $\mathcal{M}_2$ and $\mathcal{M}_3$ are the spatial versions with space in the either the mean or probability of site

unsuitability, respectively. To compare the spatial and nonspatial versions, we conduct simulations generating $n_{train} = 400$ observations from the spatial model, fitting the data to the true model and $\mathcal{M}_1$, and predicting on $n_{test} = 200$ points. The spatial region is a square on $(0, 50) \times (0, 50)$, and locations are generated randomly uniformly. Data are generated and the models fit using $a = 1$, $\sigma^2 = 1$, and $\phi = 20$, with single independent $N(0, 1)$ covariates in $\pi(\mathbf{s}_i)$ and $\mu(\mathbf{s}_i)$. We use a Normal prior centered at the truth with standard deviation 0.25 for the appropriate intercept when fitting either the true model ($\mathcal{M}_2$ or $\mathcal{M}_3$) and $\mathcal{M}_1$. When fitting $\mathcal{M}_2$ or $\mathcal{M}_3$, $\phi$ is held fixed at one-third the maximum observed distance.

We first compare $\mathcal{M}_2$ with $\mathcal{M}_1$. Simulations 2.1-2.3 in the Supplemental Information compare the parameter recovery of the two models for $\sigma_\eta^2 = \{0.5, 1, 2\}$. Table S5 of posterior means and 95% credible intervals reveal the ability of $\mathcal{M}_2$ to recover the true parameters, including the spatial variance. Fig. S8 shows that we are able to estimate the spatial surface well. In contrast, we find $\mathcal{M}_1$ often fails to capture the truth for the parameters involved in the Beta regression as $\sigma_\eta^2$ increases. We then evaluate the predictive performance of the two models by generating data under $\mathcal{M}_2$ with varying proportions of 0s, and evaluating $R^2$, AUC, and $CRPS$ for when fitting the data with $\mathcal{M}_2$ and $\mathcal{M}_1$, averaged across 50 runs (Table S6). $\mathcal{M}_2$ always outperforms in terms of AUC, $R^2$, and $CRPS_h$. Interestingly, when the data consist of a large proportion of zeros ($> 70\%$, bottom right cells in Table S6), the incorrect model actually performs better in terms of $CRPS_f$. In these scenarios, the prevalence of $Y_i > 0$ is sparse which inhibits learning the spatial surface. For a $Y_i = 0$, whether $\mu(\mathbf{s}_i)$ receives a random effect depends on the iteration's Bernoulli trial for the corresponding $Z_i$. We suspect this uncertainty leads to some model instability, and reflects the weakness of $\mathcal{M}_2$ in separating the two sources with a large number of zeros in the data.

We also generate data where the spatial effect is added to the regression for $\pi(\mathbf{s}_i)$ and compare $\mathcal{M}_3$ with $\mathcal{M}_1$ (Simulations 3.1-3.3, Supplement). We find that $\mathcal{M}_3$ and $\mathcal{M}_1$ perform very similarly in terms of parameter recovery (Table S7, S11, third column). Fig. S10 displays the estimated spatial surface, and S11 presents plots of comparisons from these simulations. Results of simulations comparing predictive performance under differing amounts of zeros are presented in Table S8. We find that $\mathcal{M}_3$ is marginally superior to the non-spatial model in all scenarios across the four metrics, except when there are a few number of zeros in the data (bottom right cell of Table S8). In these cases, only a fraction of the already low-prevalence 0s will arise due to unsuitability. Therefore, the estimate of $\pi_i$ parameter in the non-spatial model is only marginally impacted by the lack of spatial process, which serves to inflate/deflate the probability of zero.

# 5   Application to Cape Floristic Region data

## 5.1   $\mathcal{M}_1$: Results for the CFR data

We fit our proposed model to the percent cover data for two plant families, Crassulaceae and Restionaceae. We use average annual potential evaporation and the rainfall concentration index as predictors for unsuitability. These predictors capture the overall aridity of a site which should provide a hard threshold (particularly for the succulent Crassulaceae) as to whether certain species within the two families could establish. We use average annual precipitation and minimum temperature in July (austral winter) to estimate the mean of the Beta distribution. If a site is suitable, these variables should capture factors that would attenuate plant abundance through water availability (i.e., precipitation) or amount plants could grow in winter (proxied by the minimum winter temperature). All covariates were centered and scaled. To center the prior for the intercept in the degenerate probability, for various intercept values we fit the model to 130 points and computed the average $CRPS_f$ on the 50 held-out points. Additionally, at this stage, we fix $a = 1$, thus *a priori* giving equal support above and below 0. We present the models with the lowest $CRPS_f$; centering $\gamma_0$ at different values does not impact significance nor sign of the remaining coefficients.

Table 2a displays posterior means and 95% credible intervals for both plant families. For Crassulaceae, the large negative $\gamma_0$ intercept (-1.345) suggests that fewer zeros arose from unsuitability; most sites were climatically suitable. The negative coefficient for mean annual potential evaporation indicates that for fixed rainfall concentration, higher evaporation increases the suitability of a location, reasonable given that Crassulaceae are succulent plants adapted for drier conditions. The effect of rainfall concentration on unsuitability was insignificant. Higher annual mean precipitation and colder winter temperatures both tended to favor higher Crassulaceae percent cover, holding the other covariate constant. The Restionaceae typically had results opposite to those of Crassulaceae, and there was evidence that more concentrated rainfall throughout the year at a site led to higher suitability compared to sites with similar evaporation.

We hypothesize that opposing results between the two families are driven by both fire dynamics and subregion-specific climate influences. While members of the Restionaceae can inhabit a wide range of moisture availability levels, restio-dominated fynbos primarily occur on warmer north-facing slopes on drought-prone soils (Rebelo et al. 2006). The Restionaceae are also a fire-adapted family, e.g., members germinate in response to smoke (Brown et al. 1994) and have reseeding or resprouting life cycles (Wüest et al. 2016), that often "carry" fires in fynbos (Cowling and Holmes 1992). Our results agreed with these generalizations in that higher Restionaceae abundances were associated with lower rainfall and higher winter temperatures. Within the Baviaanskloof subregion, an additional dynamic of vegetation

turnover may play a role in the Restionaceae suitability patterns we observed. Grassy fynbos, i.e., fynbos dominated by C4 grasses, begins to outcompete restio-dominated fynbos in the eastern cape, particularly on drier north facing slopes. In this context, fynbos tends to occur in wetter areas (Rebelo et al. 2006). This may explain why we observed lower suitability for areas with higher potential evaporation, despite the favor of drier sites for abundance. On the other hand, Crassulaceae are likely "fire avoiders" and may have higher abundances in rockier areas that act as outcropping island microsites (Cousins et al. 2016). Higher Crassulaceae abundance associated with wetter, cooler sites may also reflect a broader signal of lower fire frequency in these areas.

## 5.2 $\mathcal{M}_2$ and $\mathcal{M}_3$: Results for a subset of the CFR data

Baviaanskloof and Cape Point are disjoint and too far apart geographically to envision a common spatial process for the random effects. So we perform spatial analysis only on Baviaanskloof as two-thirds ($n = 119$) of the total observations are located in this subregion. The analyses presented here are for Crassulaceae; similar analyses for Restionaceae are provided in the Supplement (Tables S9, S11 and Fig. S12).

In fitting the two spatial models to the Crassulaceae data in Baviaanskloof, we fixed $a = 1$ and experimented with different informative priors for the respective intercepts. Similar to Section 5.1, we ultimately select a prior which leads to the best $CRPS_f$. We began with the same covariates used to fit the non-spatial $\mathcal{M}_1$ in Section 5.1. When fitting $\mathcal{M}_2$, minimum July temperature was no longer significant for $\mu(\mathbf{s}_i)$. For $\mathcal{M}_3$, average annual evaporation was no longer significant for $\pi(\mathbf{s}_i)$. Otherwise, the ecological results and interpretations found in Section 5 still hold (Table 2b). The estimated spatial variance under $\mathcal{M}_2$ is larger than that of $\mathcal{M}_3$ (see Fig. 4 for estimates of the random effects).

## 5.3 Model comparison for the CFR data

Lastly, to compare the spatial and non-spatial models, we randomly split the data into a training set of 99 locations and a test set of 20 locations, and fit five models–BEZI, $\mathcal{M}_0$, $\mathcal{M}_1$, $\mathcal{M}_2$, and $\mathcal{M}_3$–using the covariates that were found significant for each respective model. Additionally, we fix $c$ at three different values: $\{0.1, 0.3, 0.5\}$ for all of the $\mathcal{M}$ models in order to assess whether the data are sensitive to the choice of $c$. We compare the models using the test data by evaluating Tjur's $R^2$ and AUC for all five models, and the two $CRPS$ metrics for $\mathcal{M}_1$-$\mathcal{M}_3$. We repeat this process thirty times and average the results. Table S10 shows that the model with the spatial effect in the mean, $\mathcal{M}_2$, with $a = 1$ outperforms the other four models in all the comparison metrics. This could be related to the estimated larger spatial variance under $\mathcal{M}_2$. Changing $c$ leads to minor differences in predictive performance.

The BEZI model exhibits the lowest performance in classifying between zero and positive percent covers for the data (lowest $R^2$ and AUC, Table S10). This suggests that modeling two sources of zeros is more appropriate for the data, with the inclusion of spatial dependence for $\mu_i$ yielding the best fit. Comparisons for Restionaceae are presented in Table S11.

# 6   Conclusions

Modelling zero-inflated percent cover data has been a statistical challenge for ecologists. While many approaches have been proposed, we have argued that these do not satisfactorily address the zero-inflation issue. We have developed a left-censored approach to obtain zeros for data on the unit interval, and offered a zero-inflated model for percent or proportion data which parallels such modeling for count data. The specifications enable better understanding of the incidence of absence as well as the explanatory contribution of environmental unsuitability versus absence by chance. Given that ecological data are often spatial, we have supplied spatial models to enhance explanation and enable prediction to unobserved sites with known environmental attributes. Better ecological insights emerge, particularly illustrated by our example with plant percent cover data.

Future work, in progress, considers joint zero-inflation. For example, we can examine percent cover for a pair of species at a site to allow deeper focus on biotic factors (i.e., interspecific competitive interactions) that influence site suitability and relative species abundances. Given the sum constraint at a site to at most 1, this framework anticipates negative association between pairs of species.

# Appendix

Here, we offer analytical insight into the effect of the choice of support for the latent variable, $W$ (suppressing the site subscript). By letting $W_a = -a + (1 + a)V$ where $V \sim Beta(\mu, \nu)$, $W$ has support $[-a, 1)$. Again, $a = 1$ provides matching support above and below 0. Then, $P(W_a \leq 0) = P(V \leq \frac{a}{a+1} \equiv c(a))$. Though $W_a$ is a linear transformation of $W_1$, there is no linear transformation of the Beta distribution that moves $c$ to .5, corresponding to $a = 1$. We cannot "relocate" $\mu$ for general $a$ to $\mu$ for $a = 1$. The shape of the Beta changes with $\mu$.

It is more convenient to work with $c$ since we seek the effect of changing $\mu$ for a fixed $\nu$ and $c$ is on the support of $V$. We seek to understand the behavior of $P(V \leq c|\mu, \nu)$. First, for a fixed value $p$ of this probability and also a fixed $\nu$, qualitatively, $c \uparrow \mu$. Pursuing this analytically, consider the implicit function, $F(c_{p,\nu}(\mu); \mu, \nu) = p$ where $F$ denotes a Beta cdf. Plots of the function, $c_{p,\nu}(\mu)$ vs. $\mu$, from $(0, 1)$ to $(0, 1)$ appear sigmoidal, i.e., with

asymptotes at 0 and at 1 (see Figure S13).

$S$-shaped curves are often characterized by an ordinary differential equation. For example, the generalized logistic functions (Richards' curves) are a rich class including the Gompertz and Weibull with such an attractive characterization as well as an explicit solution (a convenient current cite is $https://en.wikipedia.org/wiki/Generalisedlogisticfunction$). Here, we produce the differential equation which $c_{p,\nu}(\mu)$ satisfies. It is not an ordinary differential equation and, in general, it has no closed form solution since $F(\cdot;\mu,\nu)$ has no closed form. We have $p = F(c_{p,\nu}(\mu);\mu,\nu) = \int_0^{c_{p,\nu}(\mu)} f(x;\mu,\nu)dx$ where $f$ is the Beta$(\mu,\nu)$ density. Using Leibniz's rule, we take the derivative of both sides yielding

$$f(c_{p,\nu}(\mu);\mu,\nu)dc_{p,\nu}(\mu))/d\mu + \int_0^{c_{p,\nu}(\mu)} \frac{\partial f(x;\mu,\nu)}{\partial x}dx = 0.$$

Solving, we have the differential equation

$$dc_{p,\nu}(\mu))/d\mu = \frac{-\int_0^{c_{p,\nu}(\mu)} \frac{\partial f(x;\mu,\nu)}{\partial x}dx}{f(c_{p,\nu}(\mu);\mu,\nu)}. \tag{4}$$

Restoring the $(\alpha,\beta)$ notation, the integral in the numerator exists if $\alpha > 1$ and $\beta > 1$ (implying $\nu > 2$); otherwise $c_{p,\nu}(\mu)$ is not differentiable. If the integral exists, we have

$$-\int_0^{c_{p,\nu}(\mu)} \frac{\partial f(x;\mu,\nu)}{\partial x}dx = (\alpha+\beta+1)[F(c_{p,\nu}(\mu);\alpha,\beta-1) - F(c_{p,\nu}(\mu);\alpha-1,\beta)]. \tag{5}$$

Clearly this expression is $> 0$ and since the denominator is positive, the monotone increasing behavior of $c_{p,\nu}(\mu)$ is demonstrated. The support for the existence over $(\nu,\mu) \in (0,\infty)\times(0,1)$ satisfies $\mu > 1/\nu$ and $\mu < 1 - 1/\nu$.
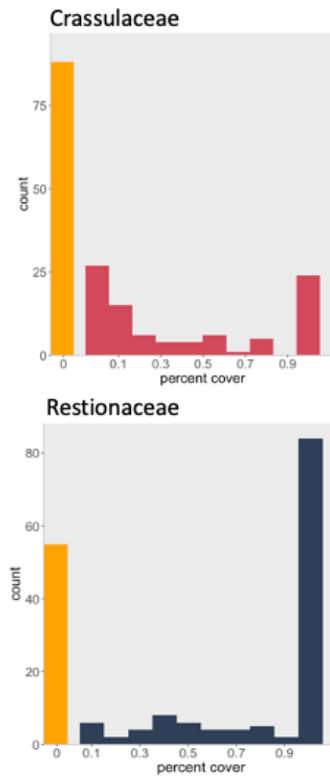
For $c_{p,\nu}(\mu)$ to be $S$-shaped, we need to show that it has exactly one point of inflection, i.e., that $dc_{p,\nu}(\mu))/d\mu$ increases up to a point and then decreases. Again using Leibniz's rule, we can take the derivative of $[F(c_{p,\nu}(\mu);\alpha,\beta-1) - F(c_{p,\nu}(\mu);\alpha-1,\beta)]$. This derivative exists if $\alpha > 2$ and $\beta > 2$. So, now $\nu > 4$ and the support for the existence of the second derivative is $\mu > 2/\nu$ and $\mu < 1 - 2/\nu$. Omitting details, in special cases we can demonstrate that this derivative starts positive and then becomes negative, implying a unique point of inflection.

Finally, to return to sensitivity to choice of $c$, we fit $S-$shaped curves to $c_{p,\nu}(\mu)$ for various $p$ and $\nu > 4$. In the literature, the support of the $S$-shaped curves is always on $R^1$ while we need support $(0,1)$. We propose to use the logit function $g_1(x) = \ln\frac{x}{1-x}$ from $(0,1)$ to $R^1$. Then, we apply the generalized logistic function, $g_2(t) = (1 + \beta_0 e^{-\beta_1 t})^{-1/\gamma}$. Plugging in, we obtain the class $g_2(g_1(x)) = \left(1 + \beta_0\left(\frac{1-x}{x}\right)^{\beta_1}\right)^{-1/\gamma}$. This three parameter function is
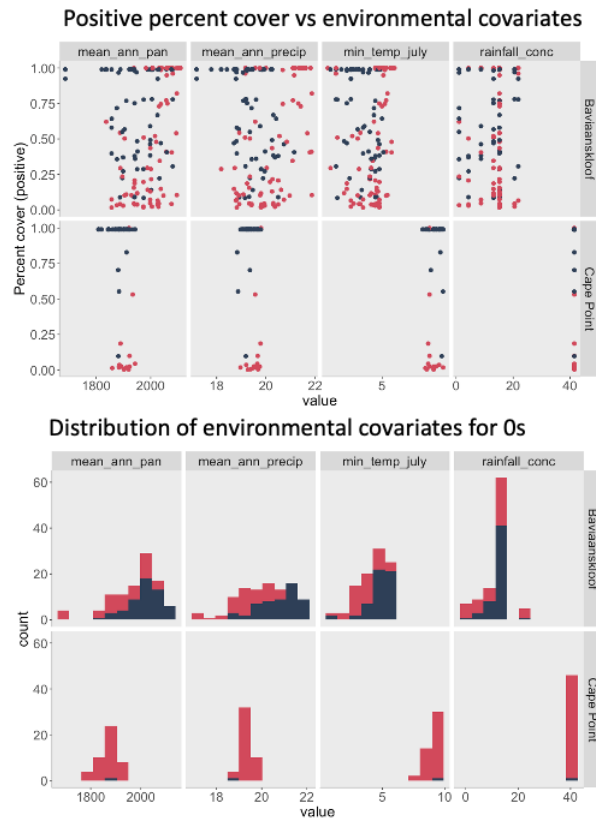
17

straightforward to fit to $c_{p,\nu}(\mu)$. Implementing the curve fitting over a range of $p$ and $\nu$ (Figure S14), we see that the curves are essentially indistinguishable implying that there is little sensitivity to choice of $a$.

Figure 2: Exploratory analysis of observed percent cover for Crassulaceae (pink) and Restionaceae (navy). Where appropriate, orange corresponds to observed 0 percent cover.

(b) Top: positive percent cover plotted against environmental covariates. Bottom: distribution of environmental covariates at locations with observed 0.

(a) Histograms of percent cover by species across the two regions.



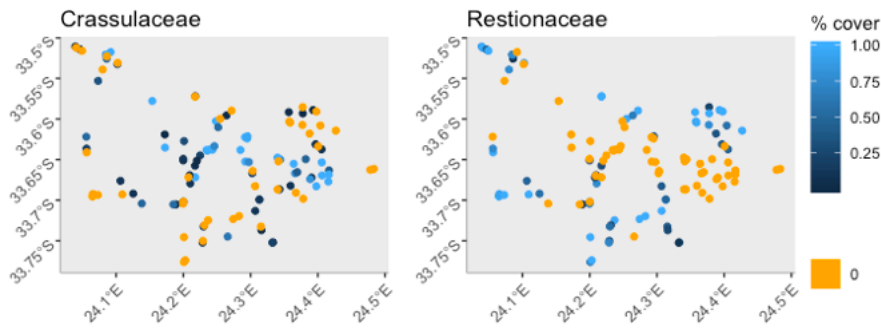(c) Spatial maps of percent cover within Baviaanskloof.



19

Figure 3: Example of the LECB model. A Beta(3,2) distribution is extended to (-1,1) using $a = 1$ and subsequently left-censored. **a)** Probability and cumulative distribution functions. **b)** Histograms of randomly generated data.

(a) PDF of Beta(3,2) (left) and PDF (middle) and CDF (right) of Extended Beta(3,2). Dark-grey shaded area denotes the left-censoring associated with the Beta mass at 0.



(b) Histograms of data generated under the Beta(3,2) distribution (left) and subsequently extended and left-censored (right).

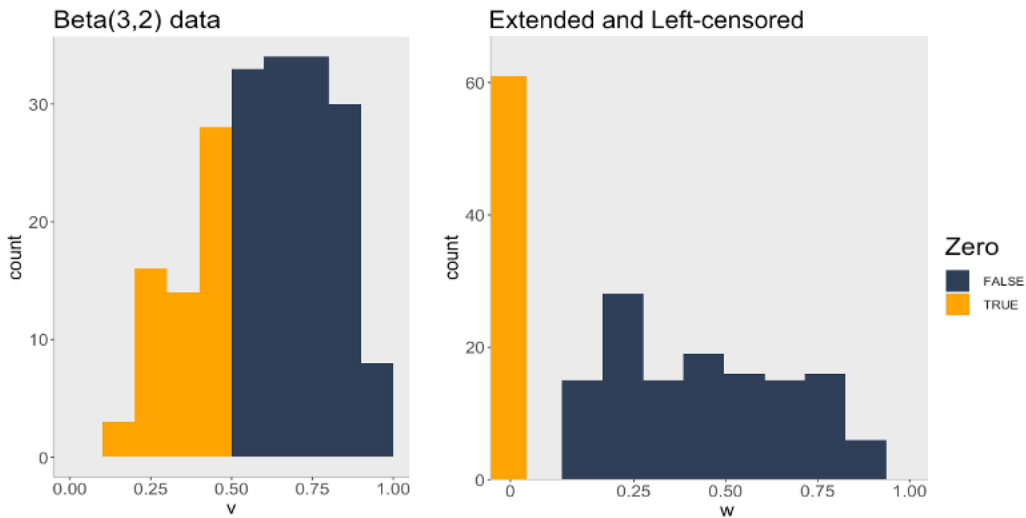Table 2: Posterior means and 95% credible intervals for coefficients from fitting various zero-inflated Beta models to CFR data. For all the following, $c$ was fixed at 0.5. The $\gamma$ coefficients are for probability of unsuitability, the $\delta$ and $\psi$ coefficients are for the Beta regression. Dashes correspond to intervals that included 0.
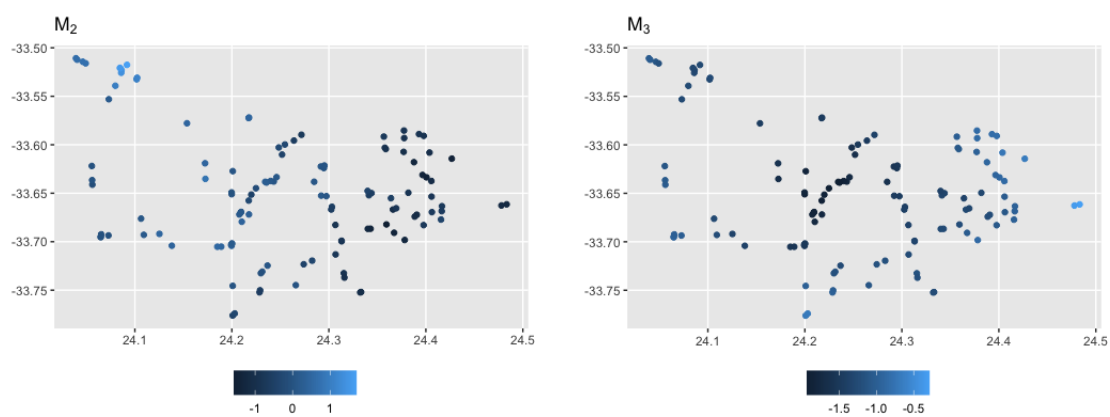
(a) Posterior summaries from non-spatial model $\mathcal{M}_1$ fit independently for each species across both Baviaanskloof and Cape Point. Priors for $\gamma_0$ for both species are centered at $-0.75$.

|  | Crassulaceae | | | Restionaceae | | |
|---|---|---|---|---|---|---|
|  | mean | 2.5% | 97.5% | mean | 2.5% | 97.5% |
| $\gamma_0$: Intercept | -1.345 | -1.765 | -0.876 | -1.0058 | -1.2453 | -0.7412 |
| $\gamma_1$: mean_ann_pan | -1.269 | -1.960 | -0.377 | 0.884 | 0.294 | 1.630 |
| $\gamma_2$: rainfall_conc | - | - | - | -0.676 | -1.210 | -0.015 |
| $\delta_0$: Intercept | 0.238 | 0.041 | 0.473 | 1.685 | 1.420 | 2.115 |
| $\delta_1$: mean_ann_precip | 0.800 | 0.625 | 0.998 | -0.896 | -1.233 | -0.505 |
| $\delta_2$: min_temp_july | -0.396 | -0.694 | -0.061 | 0.419 | 0.269 | 0.583 |
| $\psi$ | 1.023 | 0.790 | 1.306 | 1.360 | 1.036 | 1.776 |

(b) Posterior summaries from the spatial zero-inflated Beta models $\mathcal{M}_2$ and $\mathcal{M}_3$ fit to the Crassulaceae percent cover data in Baviaanskloof. $\mathcal{M}_2$ includes the spatial effect in $\mu_i$, with the prior for $\delta_0$ centered at 1. $\mathcal{M}_3$ includes the effect in $\pi_i$, with the prior for $\gamma_0$ centered at -0.75. For both models, $\phi$ fixed at one-third the maximum observed distance (14.8 km).

|  | $\mathcal{M}_2$ | | | $\mathcal{M}_3$ | | |
|---|---|---|---|---|---|---|
|  | mean | 2.5% | 97.5% | mean | 2.5% | 97.5% |
| $\gamma_0$: Intercept | -0.951 | -1.326 | -0.572 | -1.080 | -1.642 | -0.377 |
| $\gamma_1$: mean_ann_pan | -0.787 | -1.395 | -0.147 | - | - | - |
| $\delta_0$: Intercept | 0.956 | 0.658 | 1.271 | 0.574 | 0.312 | 0.855 |
| $\delta_1$: mean_ann_precip | 1.222 | 0.903 | 1.677 | 1.140 | 0.819 | 1.504 |
| $\delta_2$: min_temp_july | - | - | - | -0.405 | -0.683 | -0.056 |
| $\psi$ | 1.628 | 1.230 | 2.041 | 0.920 | 0.647 | 1.292 |
| $\sigma^2$ | 1.510 | 1.102 | 2.047 | 0.920 | 0.647 | 1.292 |

Figure 4: Posterior means of spatial random effects under the two spatial models $\mathcal{M}_2$ (left) and $\mathcal{M}_3$ (right) fit on the Crassulaceae data in Baviaanskloof. Please see Table S10

# References

Agarwal, D. K., Gelfand, A. E. and Citron-Pousty, S. (2002), 'Zero-inflated models with application to spatial count data', *Environmental and Ecological statistics* **9**(4), 341–355.

Amemiya, T. (1984), 'Tobit models: A survey', *Journal of econometrics* **24**(1-2), 3–61.

Banerjee, S., Carlin, B. P. and Gelfand, A. E. (2014), *Hierarchical modeling and analysis for spatial data*, CRC press.

Bell, J. D. and Galzin, R. (1984), 'Influence of live coral cover on coral-reef fish communities', *Marine Ecology Progress Series* **15**, 265–274.

Blasco-Moreno, A., Pérez-Casany, M., Puig, P., Morante, M. and Castells, E. (2019), 'What does a zero mean? Understanding false, random and structural zeros in ecology', *Methods in Ecology and Evolution* **10**(7), 949–959.

Born, J., Linder, H. P. and Desmet, P. (2006), 'The Greater Cape Floristic Region: Greater Cape Floristic Region', *Journal of Biogeography* **34**(1), 147–162.

Brown, N. A. C., Jamieson, H. and Botha, P. A. (1994), 'Stimulation of seed germination in South African species of Restionaceae by plant-derived smoke', *Plant Growth Regulation* **15**(1), 93–100.

Chib, S. (1992), 'Bayes inference in the tobit censored regression model', *Journal of Econometrics* **51**(1-2), 79–99.

Cousins, S. R., Witkowski, E. T. F. and Pfab, M. F. (2016), 'Beating the blaze: Fire survival in the fan aloe (*Kumara plicatilis*), a succulent monocotyledonous tree endemic to the Cape fynbos, South Africa', *Austral Ecology* **41**(5), 466–479.

Cowling, R. M. and Holmes, P. M. (1992), Flora and vegetation, *in* 'The ecology of fynbos: nutrients, fire and diversity', Oxford University Press, Cape Town, pp. 23–61.

Damgaard, C. and Irvine, K. (2019), 'Using the beta distribution to analyse plant cover data', *Journal of Ecology* **107**, 2747–2759.

Dengler, J., Jansen, F., Glöckler, F., Peet, R. K., De Cáceres, M., Chytrỳ, M., Ewald, J., Oldeland, J., Lopez-Gonzalez, G., Finckh, M. et al. (2011), 'The global index of vegetation-plot databases (givd): a new resource for vegetation science', *Journal of Vegetation Science* **22**(4), 582–597.

Douma, J. C. and Weedon, J. T. (2019), 'Analysing continuous proportions in ecology and evolution: A practical introduction to beta and Dirichlet regression', *Methods in Ecology and Evolution* **10**(9), 1412–1430.

Ellenberg, H. and Mueller-Dombois, D. (1974), Community sampling: The Relevé method, *in* 'Aims and methods of vegetation ecology', John Wiley & Sons, New York.

Ferrari, S. and Cribari-Neto, F. (2004), 'Beta regression for modelling rates and proportions', *Journal of applied statistics* **31**(7), 799–815.

Ghosh, S. K., Mukhopadhyay, P. and Lu, J.-C. J. (2006), 'Bayesian analysis of zero-inflated regression models', *Journal of Statistical planning and Inference* **136**(4), 1360–1375.

Gneiting, T. and Raftery, A. E. (2007), 'Strictly proper scoring rules, prediction, and estimation', *Journal of the American statistical Association* **102**(477), 359–378.

Hall, D. B. (2000), 'Zero-inflated poisson and binomial regression with random effects: a case study', *Biometrics* **56**(4), 1030–1039.

Hijmans, R. J., Cameron, S. E., Parra, J. L., Jones, P. G. and Jarvis, A. (2005), 'Very high resolution interpolated climate surfaces for global land areas', *International Journal of Climatology: A Journal of the Royal Meteorological Society* **25**(15), 1965–1978.

Jenkins, J. C., Chojnacky, D. C., Heath, L. S. and Birdsey, R. A. (2003), 'National-scale biomass estimators for United States tree species', *Forest science* **49**(1), 12–35.

Lambert, D. (1992), 'Zero-inflated poisson regression, with an application to defects in manufacturing', *Technometrics* **34**(1), 1–14.

Linder, H. (2001), *The African Restionaceae: An interactive key to the species (on CD ROM)*, Vol. 20.

Long, J. S. and Long, J. S. (1997), *Regression models for categorical and limited dependent variables*, Vol. 7, Sage.

Martin, T. G., Wintle, B. A., Rhodes, J. R., Kuhnert, P. M., Field, S. A., Low-Choy, S. J., Tyre, A. J. and Possingham, H. P. (2005), 'Zero tolerance ecology: improving ecological inference by modelling the source of zero observations: Modelling excess zeros in ecology', *Ecology Letters* **8**(11), 1235–1246.

Mullahy, J. (1986), 'Specification and testing of some modified count data models', *Journal of econometrics* **33**(3), 341–365.

Murray, I., Prescott Adams, R. and MacKay, D. J. (2010), 'Elliptical slice sampling'.

Ospina, R. and Ferrari, S. L. (2010), 'Inflated beta distributions', *Statistical papers* **51**(1), 111.

Ospina, R. and Ferrari, S. L. (2012), 'A general class of zero-or-one inflated beta regression models', *Computational Statistics I& Data Analysis* **56**(6), 1609–1623.

Peterson, E. B. (2005), 'Estimating cover of an invasive grass (*Bromus tectorum* using tobit regression and phenology derived from two dates of Landsat ETM+ data', *International Journal of Remote Sensing* **26**(12), 2491–2507.

R Core Team (2013), *R: A language and environment for statistical computing*, R Foundation for Statistical Computing, Vienna, Austria.

Rathbun, S. L. and Fei, S. (2006), 'A spatial zero-inflated poisson regression model for oak regeneration', *Environmental and Ecological Statistics* **13**(4), 409–426.

Rebelo, A. G., Boucher, C., Helme, N., Mucina, L. and Rutherford, M. C. (2006), Fynbos Biome, *in* 'The vegetation of South Africa, Lesotho, and Swaziland', South African National Biodiversity Institue, Pretoria, South Africa.

Schaminée, J. H., Hennekens, S. M., Chytry, M. and Rodwell, J. S. (2009), 'Vegetation-plot data and databases in europe: an overview'.

Schulze, R. E. (1997), South African atlas of agrohydrology and climatology: Contribution towards a final report to the water research commission on project 492, Technical Report TT82-96, Water Resource Commission, Pretoria, South Africa.

Siegfried, S. and Hothorn, T. (2020), 'Count transformation models', *Methods in Ecology and Evolution* **11**(7), 818–827.

Tjur, T. (2009), 'Coefficients of determination in logistic regression models—a new proposal: The coefficient of discrimination', *The American Statistician* **63**(4), 366–372.

Tobin, J. (1958), 'Estimation of relationships for limited dependent variables', *Econometrica: journal of the Econometric Society* pp. 24–36.

Turpie, J. K., Heydenrych, B. J. and Lamberth, S. J. (2003), 'Economic value of terrestrial and marine biodiversity in the cape floristic region: implications for defining effective and socially optimal conservation strategies', *Biological conservation* **112**(1-2), 233–251.

van der Maarel, E. (2007), 'Transformation of cover-abundance values for appropriate numerical treatment- Alternatives to the proposals by Podani', *Journal of Vegetation Science* **18**, 767–770.

Vanha-Majamaa, I., Salemaa, M., Tuominen, S. and Mikkola, K. (2000), 'Digitized photographs in vegetation analysis- a comparison of cover estimates', *Applied Vegetation Science* **3**, 89–94.

Veech, J. A., Ott, J. R. and Troy, J. R. (2016), 'Intrinsic heterogeneity in detection probability and its effect on $N$-mixture models', *Methods in Ecology and Evolution* **7**(9), 1019–1028.

Wenger, S. J. and Freeman, M. C. (2008), 'Estimating species occurrence, abundance, and detection probability using zero-inflated distributions', *Ecology* **89**(10), 2953–2959. Publisher: Wiley Online Library.

Wüest, R. O., Litsios, G., Forest, F., Lexer, C., Linder, H. P., Salamin, N., Zimmermann, N. E. and Pearman, P. B. (2016), 'Resprouter fraction in Cape Restionaceae assemblages varies with climate and soil type', *Functional Ecology* **30**(9), 1583–1592.

Xie, Y., Civco, D. L. and Silander, J. A. (2018), 'Species-specific spring and autumn leaf phenology captured by time-lapse digital cameras', *Ecosphere* **9**(1).
**URL:** *https://onlinelibrary.wiley.com/doi/10.1002/ecs2.2089*

Xie, Y., Wang, X., Wilson, A. M. and Silander, J. A. (2018), 'Predicting autumn phenology: How deciduous tree species respond to weather stressors', *Agricultural and Forest Meteorology* **250-251**, 127–137.
**URL:** *https://linkinghub.elsevier.com/retrieve/pii/S0168192317306792*

Zhang, H. (2004), 'Inconsistent estimation and asymptotically equal interpolations in model-based geostatistics', *Journal of the American Statistical Association* **99**(465), 250–261.

# Supplementary Information

Becky Tang[1*], Henry A. Frye[2], Alan E. Gelfand[1], John A Silander, Jr.[2]

## 1  Monte Carlo integration for CRPS

For observation $y_i$ and given $N$ uniform samples $x^{(j)}$,

$$\int_0^\infty (F(x) - \mathbf{1}(y < x))^2 dx \approx \frac{1}{N} \sum_{j=1}^N \left( F(x^{(j)}; \alpha_i, \beta_i) - \mathbf{1}(y_i < x^{(j)}) \right)^2 \tag{1}$$

If the posterior predictive CDF $F$ does not have a closed form, then given a sequence of $\{\theta_j\}_{j=1}^m$ of parameter values from the posterior distribution, $F$ can be approximated via $\hat{F}(x) \approx \frac{1}{m} \sum_{j=1}^m F_c(x|\theta_j)$, where $F_c(x|\theta)$ is the conditional predictive CDF.

# 2 Simulation details for comparing BEZI, LCEB, and zero-inflated Beta models

In Section 4.2 of the main text, we present results of simulations to examine the prediction performances of the BEZI model of Ospina and Ferrari (2010), our censor-only model LCEB, and our zero-inflated Beta model $\mathcal{M}_1$. In all cases, we generate $n_{test} = 100$ and $n_{train} = 200$ observations from one model, fit all three models to the train data, and obtain predictions for the test data. We compare the predictions to the true data by calculating $R^2$ and AUC for the probabilities of classifying on observation as 0 or positive, and average $CRPS_h$ for all positive test observations. For all simulations, we set the true $\nu = 4.5$, i.e. just an intercept for the regression for $\nu$ (see the appendix for this choice of $\nu$). Chains were run for 5000 iterations after 2500 burnin, and thinned to retain every fifth sample. For models LCEB and $\mathcal{M}_1$, the data are generated and models fit with $a = 1$. For all the parameters in all three models, we use diffuse $N(0, 10)$ priors. The only exception is that we employ an informative prior for $\gamma_0$ when fitting $\mathcal{M}_1$, as described below.

For data generated from the zero-inflated Beta model $\mathcal{M}_1$, the regressions for $\pi_i$ and $\mu_i$ each contain an intercept and one independent $N(0, 1)$ covariate, i.e. $Z_i = (1 \; Z_{i1})'$ and $X_i = (1 \; X_{i1})'$, where $Z_{i1} \sim N(0, 1)$ independent of $X_{i1} \sim N(0, 1)$. The proportions of 0s and their sources are controlled by varying the intercepts $\gamma_0$ and $\delta_0$. The BEZI model is fit with the same $\mathbf{Z}$ and $\mathbf{X}$ for the regressions for $\pi$ and $\mu$, respectively. LCEB is fit by concatenating or joining the two covariate sets together, i.e. $X_{i,LCEB} = (1 \; X_{i1} \; Z_{i1})'$. When fitting $\mathcal{M}_1$, we use a $N(\gamma_{0,true}, \sigma^2 = 0.25)$ prior for $\gamma_0$.

For data generated from the BEZI model, the regressions for $\pi_i$ and $\mu_i$ once again each contain one independent $N(0, 1)$ covariate. The proportion of 0s is controlled by varying the intercept $\gamma_0$, and the magnitude of the positive percent covers is controlled by $\delta_0$. $\mathcal{M}_1$ is fit with the same $\mathbf{Z}$ and $\mathbf{X}$ for the regressions for $\pi$ and $\mu$, respectively. LCEB is fit by concatenating or joining the two covariate sets together, i.e. $X_{i,LCEB} = (1 \; X_{i1} \; Z_{i1})'$. When fitting $\mathcal{M}_1$, we use a $N(\gamma_{0,true}, \sigma^2 = 0.25)$ prior for $\gamma_0$.

For data generated from the censor-only model LCEB, we only have a single covariate set $\mathbf{X}$, where $X_i = (1 \; X_{i1})'$ and $X_{i1} \sim N(0, 1)$. The BEZI and $\mathcal{M}_1$ models are fit with $\mathbf{Z} = \mathbf{X}$. When fitting $\mathcal{M}_1$, we use a $N(0, \sigma^2 = 0.25)$ prior for $\gamma_0$.

# 3    Posterior distributions for spatial zero-inflated Beta

If we place the spatial random effects $\boldsymbol{\eta}(\mathbf{s})$ into the mean $\mu(\mathbf{s}_i)$ and hold the spatial decay parameter $\phi_\eta$ fixed, then the posterior distribution is proportional to:

$$\prod_{i=1}^{n} f_{Y_i|Z_i,\boldsymbol{\delta},\boldsymbol{\psi}}(Y_i|Z_i,\boldsymbol{\delta},\boldsymbol{\psi},\eta(\mathbf{s}_i)) \cdot \prod_{i=1}^{n} f_{Z_i|\boldsymbol{\gamma}}(Z_i|\boldsymbol{\gamma}) \cdot f_{\boldsymbol{\eta}|\sigma^2}(\boldsymbol{\eta}(\mathbf{s})|\sigma_\eta^2) \cdot f_{\boldsymbol{\gamma}}(\boldsymbol{\gamma}) \cdot f_{\boldsymbol{\delta}}(\boldsymbol{\delta}) \cdot f_{\boldsymbol{\psi}}(\boldsymbol{\psi}) \cdot f_{\sigma^2}(\sigma_\eta^2).$$

$$(2)$$

We could also use a Gaussian process for the spatial effects $\boldsymbol{\tau}(\mathbf{s})$ in $\pi(\mathbf{s}_i)$. Once again holding the spatial decay parameter $\phi_\tau$ fixed, the posterior is proportional to:

$$\prod_{i=1}^{n} f_{Y_i|Z_i,\boldsymbol{\delta},\boldsymbol{\psi}}(Y_i|Z_i,\boldsymbol{\delta},\boldsymbol{\psi}) \cdot \prod_{i=1}^{n} f_{Z_i|\boldsymbol{\gamma},\tau(\mathbf{s}_i)}(Z_i|\boldsymbol{\gamma},\tau(\mathbf{s}_i)) \cdot f_{\boldsymbol{\tau}|\sigma^2}(\boldsymbol{\tau}(\mathbf{s})|\sigma_\tau^2) \cdot f_{\boldsymbol{\gamma}}(\boldsymbol{\gamma}) \cdot f_{\boldsymbol{\delta}}(\boldsymbol{\delta}) \cdot f_{\boldsymbol{\psi}}(\boldsymbol{\psi}) \cdot f_{\sigma^2}(\sigma_\tau^2).$$

$$(3)$$

# 4 Supplementary figures

Figure 1: Histograms of environmental covariates within Baviaanskloof (grey) and Cape Point (green).
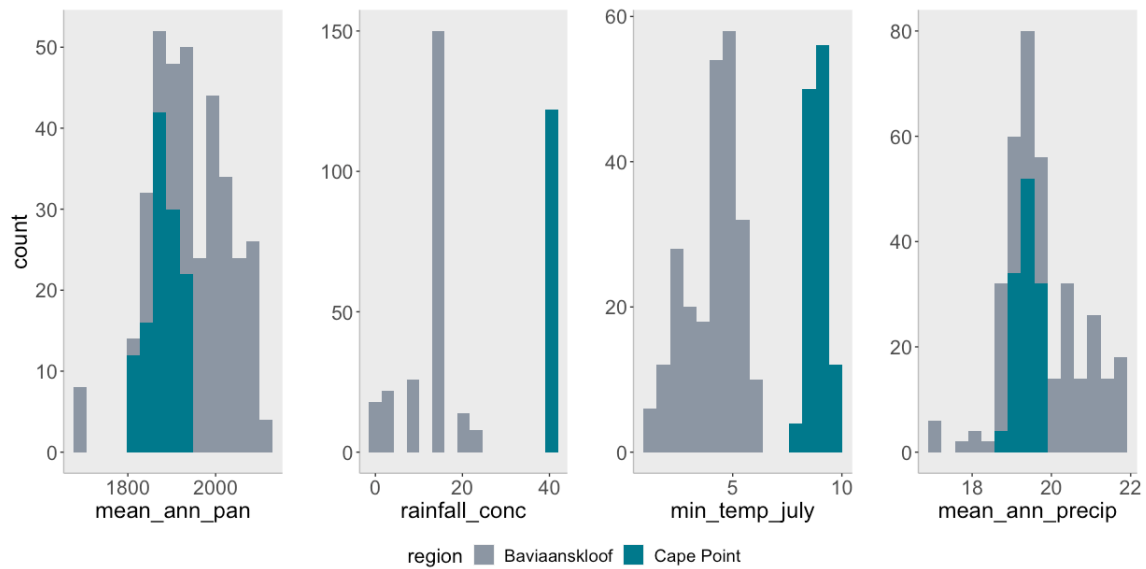
Figure 2: Spatial plots of percent cover by species in Cape Point. Orange points denote observed 0, and blue points are positive covers.

Figure 3: Results of simulation study for censored model. Data generated under various values of $c$ (panel titles), then fit using models with $c$ held fixed (x-axis). Regression for $\mu$ has a single covariate, i.e. $logit(\mu) = \delta_0 + \delta_1 x_1$, with varying $\delta_0$ and $\psi$. the average proportion of observed 0s in the training data. For Simulations 1-6: $\delta_0 = \{-0.5, 0.5, 1.5, -0.5, 0.5, 1.5\}$ and $\phi = \{3, 3, 3, 15, 15, 15\}$. Proportions of 0 in data increase with increasing true $c$.

(a) Empirical credible interval coverage for $\delta_1 = -0.5$ with nominal 95% coverage (dashed line). In all simulations, empirical coverage is closest to nominal when $c_{fixed}$ is set equal to the true $c$, as expected.



(b) Average RMSPE for the probability of zero for test data. Prediction performance tends to worsen as the magnitude $|c_{fixed} - c_{true}|$ increases (first and last panels).

Figure 4: Plots display the proportion of simulations where sign of $\delta_1 = -0.5$ was correctly identified according to 95% credible intervals. Data were generated under various values of $c$ (panel titles), then fit using models with $c$ held fixed (x-axis) with $logit(\mu) = \delta_0 + \delta_1 x_1$.

Figure 5: Predicted and true probabilities of 0 for data generated in Simulation 1. Each plot corresponds to data generated with different $c_{\text{true}}$, with colors denoting posterior mean probabilities for models fit with different $c_{\text{fixed}}$. For a given $c_{\text{true}}$, predicted probabilities are generally quite similar across the $c_{\text{fixed}}$. The notable exception is the under-prediction from the model fit with $c_{\text{fixed}} = 0.1$ to data generated with $c_{\text{true}} = 0.9$.

Figure 6: ROC curves and AUC for Simulations 1.1-1.3 for determining the source of a zero. All curves lie well about the 45-degree line, with AUCs of 0.80 or higher.

Figure 7: Run times (seconds) of zero-inflated Beta models for different sizes $n$ of dataset. $\mathcal{M}_1$ refers to the non-spatial zero-inflated Beta model, whereas $\mathcal{M}_2$ has space in $\mu(\mathbf{s}_i)$ and $\mathcal{M}_3$ has space in $\pi(\mathbf{s}_i)$. Run times are averaged across data with different proportions and sources of 0s. Simulations are run using a server node with one core Intel(R) Xeon(R) CPU E5-2680 v3 @ 2.50GHz.

Figure 8: Simulations 2.1-2.3: the first two columns show the true and posterior mean predictions of the spatial random effects, using a common color scale. The third column shows the differences between the true and predicted values. That the majority of the points in the differences plots are white or pale demonstrates that we are able to recover the spatial surface well. See Table S6 for corresponding posterior summaries.

Figure 9: Simulations 2.1-2.3: red denotes to the data-generating model $\mathcal{M}_2$, and blue denotes the incorrect model $\mathcal{M}_1$. Left: ROC curves for predicting source of zero. Middle: parallel histograms reveal $\mathcal{M}_2$'s superior performance in separating zero and positive observations. Right: Plots of estimates means and probabilities of zero reveal that $\mathcal{M}_2$ better approximates the truth when compared to $\mathcal{M}_1$.

Figure 10: Simulations 3.1-3.3: the first two columns show the true and posterior mean predictions of the spatial random effects, using a common color scale. The third column shows the differences between the true and predicted values. See Table S8 for corresponding posterior summaries.

Figure 11: Simulations 3.1-3.3: green denotes to the data-generating model $\mathcal{M}_3$, and blue denotes the incorrect model $\mathcal{M}_1$. Left: ROC curves for predicting source of zero. Middle: parallel histograms reveal $\mathcal{M}_3$'s superior performance in separating zero and positive observations. Right: Plots of estimates means and probabilities of zero show that $\mathcal{M}_3$ better approximates the true probability of 0 when compared to $\mathcal{M}_1$, but both models perform similarly in recovering the true mean.

Figure 12: Posterior means of spatial random effects under the two spatial zero-inflated Beta models $\mathcal{M}_2$ (left) and $\mathcal{M}_3$ (right) fit on the Restionaceae data in Baviaanskloof.



Figure 13: Plots of functions $\mu$ vs $c_{nu,p}(\mu)$ for various $\nu$ (colors) and $p$ (panels).



41

Figure 14: Plots of $c_{\nu,p}(\mu)$ vs. $\mu$ (in blue, similar to Fig. S13), and approximations obtained using the class of curves in Eq. 6 in main text (red), for various $\nu$ (columns) and $p$ (rows). The curves are nearly indistinguishable, implying little sensitivity to the choice of $c$ (equivalently, $a$).

# 5 Supplementary tables

Table 1: Our choices of prior distributions for this paper. Here, $g_0$, $\delta_0$, and $b$ vary depending on the context (simulation vs application). The notation for the Normal distribution denotes mean and standard deviation. Dashes denote parameters that are not applicable for the given model. Note for all $\mathcal{M}$ models, $a$ and (for the spatial models, $\phi$) are assumed to be held fixed.

| Parameter | BEZI | LCEB | $\mathcal{M}_1$ | $\mathcal{M}_2$ | $\mathcal{M}_3$ |
|---|---|---|---|---|---|
| $\gamma_0$ | $N(0,10)$ | - | $N(g_0, 0.25)$ | $N(0,10)$ | $N(g_0, 0.25)$ |
| $\gamma_j, j \neq 0$ | $N(0,10)$ | - | $N(0,10)$ | $N(0,10)$ | $N(0,10)$ |
| $\delta_0$ | $N(0,10)$ | $N(0,10)$ | $N(0,10)$ | $N(d_0, 0.25)$ | $N(0,10)$ |
| $\delta_k, k \neq 0$ | $N(0,10)$ | $N(0,10)$ | $N(0,10)$ | $N(0,10)$ | $N(0,10)$ |
| $\psi$ | $N(0,10)$ | $N(0,10)$ | $N(0,10)$ | $N(0,10)$ | $N(0,10)$ |
| $\sigma^2$ | - | - | - | $Inv.Ga(0.5, 0.5)$ | $Inv.Ga(0.5, 0.5)$ |

Table 2: Posterior mean and 95% credible intervals for parameters in Simulations 1.1-1.3.

| | Simulation 1.1 | | | | Simulation 1.2 | | | | Simulation 1.3 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 27% unsuitable, 27% ecological | | | | 12% unsuitable, 29% ecological | | | | 40% unsuitable, 10% ecological | | | |
| | true | mean | 2.5% | 97.5% | true | mean | 2.5% | 97.5% | true | mean | 2.5% | 97.5% |
| $\gamma_0$ | -1.25 | -1.280 | -1.606 | -0.936 | -2.50 | -2.526 | -2.978 | -2.097 | -0.50 | -0.604 | -0.834 | -0.351 |
| $\gamma_1$ | 0.75 | 0.665 | 0.447 | 0.977 | 0.75 | 0.705 | 0.313 | 1.139 | 0.75 | 0.716 | 0.534 | 0.935 |
| $\delta_0$ | 0.50 | 0.443 | 0.326 | 0.590 | 0.50 | 0.463 | 0.335 | 0.596 | 1.25 | 1.255 | 1.099 | 1.404 |
| $\delta_1$ | -0.50 | -0.481 | -0.558 | -0.395 | -0.50 | -0.487 | -0.558 | -0.410 | -0.50 | -0.495 | -0.579 | -0.416 |
| $\psi$ | 1.50 | 1.547 | 1.231 | 1.835 | 1.50 | 1.579 | 1.374 | 1.776 | 1.50 | 1.589 | 1.368 | 1.849 |

Table 3: Simulation for data generated under the BEZI model. $R^2$ and AUC for classifying between zero and positive, and $CRPS_h$ for predictive distribution of positive observations for various simulations. Simulations vary by amount of zeros in the data, as well the value of $\delta_0$.

| | $R^2$ | | | AUC | | | $CRPS_h$ | | |
|---|---|---|---|---|---|---|---|---|---|
| (% 0, $\delta_0$) | BEZI | LCEB | $\mathcal{M}_1$ | BEZI | LCEB | $\mathcal{M}_1$ | BEZI | LCEB | $\mathcal{M}_1$ |
| 25%, 1 | 0.023 | 0.002 | **0.024** | 0.598 | 0.562 | **0.599** | 0.076 | 0.106 | **0.053** |
| 25%, -1 | **0.031** | 0.001 | 0.025 | **0.622** | 0.547 | 0.575 | **0.054** | 0.107 | 0.057 |
| 50%, 1 | **0.041** | 0.000 | **0.041** | **0.613** | 0.554 | 0.611 | **0.053** | 0.149 | 0.107 |
| 50%, -1 | **0.035** | -0.002 | 0.030 | **0.599** | 0.533 | 0.574 | **0.054** | 0.107 | 0.073 |
| 75%, 1 | **0.030** | -0.001 | **0.030** | **0.614** | 0.549 | **0.614** | 0.110 | 0.249 | **0.054** |
| 75%, -1 | **0.023** | -0.002 | 0.021 | **0.590** | 0.563 | 0.578 | **0.055** | 0.109 | 0.065 |

Table 4: Simulation for data generated under the censor-only model LCEB. $R^2$ and AUC for classifying between zero and positive, and $CRPS_h$ for predictive distribution of positive observations for various simulations. Simulations vary by amount of zeros in the data.

| | $R^2$ | | | AUC | | | $CRPS_h$ | | |
|---|---|---|---|---|---|---|---|---|---|
| % 0 | BEZI | LCEB | $\mathcal{M}_1$ | BEZI | LCEB | $\mathcal{M}_1$ | BEZI | LCEB | $\mathcal{M}_1$ |
| 25% | 0.120 | 0.120 | **0.122** | **0.722** | **0.722** | **0.722** | 0.085 | 0.083 | **0.075** |
| 50% | 0.164 | 0.164 | **0.165** | **0.736** | **0.736** | **0.736** | 0.140 | 0.091 | **0.072** |
| 75% | 0.142 | **0.135** | 0.138 | **0.735** | **0.735** | **0.735** | 0.130 | 0.104 | **0.085** |

Table 5: Posterior summaries for coefficients for simulations where data are generated under $\mathcal{M}_2$, with simulations differing in zero amounts and sources. Simulation 2.1: 20.3% of data are unsuitable zeros and 28.3% censored zeros. Simulation 2.2: 30.7% of data are unsuitable zeros and 11% censored zeros. Simulation 2.3: 29.5% of data are unsuitable zeros and 11% censored zeros.

| | | | $\mathcal{M}_2$ | | | $\mathcal{M}_1$ | | |
|---|---|---|---|---|---|---|---|---|
| | | true | mean | 2.5% | 97.5% | mean | 2.5% | 97.5% |
| Simulation 2.1 | $\gamma_0$ | -1.000 | -1.275 | -1.683 | -0.851 | -1.121 | -1.618 | -0.670 |
| | $\gamma_1$ | 0.750 | 0.579 | 0.265 | 1.004 | 0.588 | 0.263 | 1.046 |
| | $\delta_0$ | 1.500 | 1.427 | 1.093 | 1.925 | 0.733 | 0.534 | 0.952 |
| | $\delta_1$ | -1.000 | -0.959 | -1.103 | -0.796 | -0.851 | -1.020 | -0.676 |
| | $\psi$ | 2.000 | 1.926 | 1.664 | 2.187 | 0.930 | 0.707 | 1.156 |
| | $\sigma_\eta^2$ | 1.000 | 1.016 | 0.782 | 1.352 | - | - | - |
| | $\phi_\eta$ | 15.000 | 19.830 | 19.830 | 19.830 | - | - | - |
| Simulation 2.2 | $\gamma_0$ | -1.000 | -1.031 | -1.345 | -0.704 | -1.460 | -1.992 | -0.854 |
| | $\gamma_1$ | 0.750 | 1.103 | 0.724 | 1.538 | 1.364 | 0.928 | 2.043 |
| | $\delta_0$ | 1.500 | 1.431 | 1.136 | 1.820 | 1.187 | 0.993 | 1.409 |
| | $\delta_1$ | -1.000 | -1.015 | -1.151 | -0.878 | -0.920 | -1.051 | -0.766 |
| | $\psi$ | 2.000 | 2.001 | 1.780 | 2.258 | 1.219 | 0.964 | 1.515 |
| | $\sigma_\eta^2$ | 0.500 | 0.496 | 0.283 | 0.799 | - | - | - |
| | $\phi_\eta$ | 15.000 | 20.071 | 20.071 | 20.071 | - | - | - |
| Simulation 2.3 | $\gamma_0$ | -2.000 | -1.374 | -2.110 | -0.629 | -0.298 | -0.745 | 0.174 |
| | $\gamma_1$ | 0.750 | 0.475 | 0.031 | 1.172 | 0.247 | -0.038 | 0.624 |
| | $\delta_0$ | 1.500 | 1.367 | 1.006 | 1.796 | 0.291 | 0.056 | 0.553 |
| | $\delta_1$ | -1.000 | -0.995 | -1.151 | -0.831 | -0.821 | -0.980 | -0.637 |
| | $\psi$ | 2.000 | 1.927 | 1.467 | 2.498 | 1.337 | 1.085 | 1.655 |
| | $\sigma_\eta^2$ | 2.000 | 2.010 | 1.451 | 2.704 | - | - | - |
| | $\phi_\eta$ | 15.000 | 20.311 | 20.311 | 20.311 | - | - | - |

Table 6: Metrics for comparing models $\mathcal{M}_1$ and $\mathcal{M}_2$ for data generated under $\mathcal{M}_2$ with $\sigma_\eta^2 = 1$ and $\phi_\eta = 15$ in $(0, 50) \times (0, 50)$ region. Comparison metrics of area under the ROC curve (AUC) for determining source of zero, Tjur's $R^2$ for classifying between zero and positive, and $CRPS_h$ and $CRPS_f$ for predictive distributions. Simulations vary by amount of zeros in the data, as well as the proportions of zeros arising from each source (% unsuitable and % unsuitable zeros), controlled largely by varying intercepts $\gamma_0$ and $\delta_0$. Larger AUC and $R^2$ and smaller $CRPS$ indicate superior performance (in bold).

| | AUC | | $R^2$ | | $CRPS_h$ | | $CRPS_f$ | |
|---|---|---|---|---|---|---|---|---|
| (% deg. 0, % Beta 0) | $\mathcal{M}_1$ | $\mathcal{M}_2$ | $\mathcal{M}_1$ | $\mathcal{M}_2$ | $\mathcal{M}_1$ | $\mathcal{M}_2$ | $\mathcal{M}_1$ | $\mathcal{M}_2$ |
| (14%, 9%) | 0.815 | **0.877** | 0.116 | **0.173** | 0.064 | **0.053** | 0.233 | **0.228** |
| (9%, 10%) | 0.835 | **0.893** | 0.116 | **0.192** | 0.064 | **0.053** | 0.213 | **0.207** |
| (6%, 13%) | 0.807 | **0.878** | 0.118 | **0.201** | 0.060 | **0.051** | 0.188 | **0.181** |
| (33%, 16%) | 0.765 | **0.822** | 0.106 | **0.156** | 0.078 | **0.065** | 0.317 | **0.316** |
| (21%, 23%) | 0.786 | **0.830** | 0.151 | **0.236** | 0.083 | **0.067** | 0.347 | **0.345** |
| (15%, 33%) | 0.719 | **0.768** | 0.188 | **0.294** | 0.093 | **0.075** | 0.462 | **0.460** |
| (66%, 9%) | 0.767 | **0.804** | 0.068 | **0.084** | 0.181 | **0.081** | **0.224** | 0.227 |
| (40%, 33%) | 0.721 | **0.754** | 0.126 | **0.174** | 0.093 | **0.079** | **0.364** | 0.370 |
| (16%, 60%) | 0.632 | **0.643** | 0.200 | **0.272** | 0.111 | **0.096** | **0.851** | 0.870 |

Table 7: Posterior summaries for coefficients for simulations where data are generated under $\mathcal{M}_3$, with differing in zero amounts and sources. Simulation 3.1: 34% of data are unsuitable zeros and 11.3% censored zeros. Simulation 3.2: 12% of data are unsuitable zeros and 25.3% censored zeros. Simulation 3.3: 22.3% of data are unsuitable zeros and 26.7% censored zeros.

| | | | $\mathcal{M}_3$ | | | $\mathcal{M}_1$ | | |
|---|---|---|---|---|---|---|---|---|
| | | true | mean | 2.5% | 97.5% | mean | 2.5% | 97.5% |
| Simulation 3.1 | $\gamma_0$ | -1.000 | -0.972 | -1.347 | -0.507 | -0.863 | -1.108 | -0.592 |
| | $\gamma_1$ | 0.750 | 0.736 | 0.427 | 1.147 | 0.715 | 0.458 | 1.077 |
| | $\delta_0$ | 1.500 | 1.414 | 1.246 | 1.593 | 1.351 | 1.185 | 1.523 |
| | $\delta_1$ | -1.000 | -0.987 | -1.085 | -0.877 | -0.984 | -1.077 | -0.866 |
| | $\psi$ | 2.000 | 2.075 | 1.837 | 2.324 | 1.978 | 1.714 | 2.303 |
| | $\sigma_\tau^2$ | 1.000 | 0.930 | 0.422 | 1.520 | - | - | - |
| | $\phi_\tau$ | 15.000 | 20.311 | 20.311 | 20.311 | - | - | - |
| Simulation 3.2 | $\gamma_0$ | -2.000 | -2.006 | -2.391 | -1.574 | -1.909 | -2.197 | -1.559 |
| | $\gamma_1$ | 0.750 | 0.255 | -0.216 | 0.838 | 0.178 | -0.256 | 0.742 |
| | $\delta_0$ | 1.000 | 1.019 | 0.874 | 1.176 | 1.009 | 0.871 | 1.156 |
| | $\delta_1$ | -1.000 | -0.962 | -1.045 | -0.873 | -0.969 | -1.050 | -0.865 |
| | $\psi$ | 2.000 | 2.023 | 1.793 | 2.291 | 2.020 | 1.797 | 2.266 |
| | $\sigma_\tau^2$ | 0.500 | 0.497 | 0.141 | 1.172 | - | - | - |
| | $\phi_\tau$ | 15.000 | 20.311 | 20.311 | 20.311 | - | - | - |
| Simulation 3.3 | $\gamma_0$ | -2.000 | -1.949 | -2.366 | -1.460 | -1.661 | -1.960 | -1.280 |
| | $\gamma_1$ | 0.750 | 0.911 | 0.483 | 1.452 | 0.653 | 0.218 | 1.140 |
| | $\delta_0$ | 1.000 | 0.895 | 0.788 | 1.017 | 0.819 | 0.692 | 0.987 |
| | $\delta_1$ | -1.000 | -1.090 | -1.189 | -0.973 | -1.080 | -1.165 | -0.989 |
| | $\psi$ | 2.000 | 2.164 | 1.969 | 2.402 | 2.058 | 1.838 | 2.354 |
| | $\sigma_\tau^2$ | 2.000 | 4.215 | 1.512 | 6.643 | - | - | - |
| | $\phi_\tau$ | 15.000 | 20.311 | 20.311 | 20.311 | - | - | - |

Table 8: Metrics for comparing models $\mathcal{M}_1$ and $\mathcal{M}_3$ for data generated under $\mathcal{M}_3$, with $\sigma_\tau^2 = 1$ and $\phi_\tau = 15$ in a $(0, 50) \times (0, 50)$ region. Comparison metrics of area under the ROC curve (AUC) for determining source of zero, Tjur's $R^2$ for classifying between zero and positive, and $CRPS_h$ and $CRPS_f$ for predictive distributions. Simulations vary by amount of zeros in the data as well as the proportions of zeros arising from each source (% unsuitable and % unsuitable zeros), largely controlled by varying the intercepts $\gamma_0$ and $\delta_0$. Larger AUC and $R^2$ and smaller $CRPS$ indicate superior performance (in bold).

| | AUC | | $R^2$ | | $CRPS_h$ | | $CRPS_f$ | |
|---|---|---|---|---|---|---|---|---|
| (% deg. 0, % Beta 0) | $\mathcal{M}_1$ | $\mathcal{M}_3$ | $\mathcal{M}_1$ | $\mathcal{M}_3$ | $\mathcal{M}_1$ | $\mathcal{M}_3$ | $\mathcal{M}_1$ | $\mathcal{M}_3$ |
| (20%, 8%) | 0.908 | **0.914** | 0.176 | **0.202** | 0.049 | **0.047** | 0.242 | **0.241** |
| (11%, 10%) | 0.916 | **0.919** | 0.233 | **0.244** | **0.048** | **0.048** | **0.216** | **0.216** |
| (6%, 13%) | 0.915 | **0.916** | 0.278 | **0.283** | **0.049** | **0.049** | **0.197** | **0.197** |
| (36%, 14%) | 0.857 | **0.871** | 0.181 | **0.222** | 0.058 | **0.056** | 0.300 | **0.295** |
| (23%, 24%) | 0.800 | **0.821** | 0.316 | **0.327** | 0.065 | **0.062** | 0.380 | **0.379** |
| (20%, 33%) | 0.792 | **0.813** | 0.350 | **0.360** | 0.065 | **0.064** | 0.432 | **0.422** |
| (65%, 9%) | 0.838 | **0.854** | 0.093 | **0.136** | 0.061 | **0.059** | 0.243 | **0.233** |
| (49%, 25%) | 0.772 | **0.792** | 0.249 | **0.269** | 0.067 | **0.066** | 0.382 | **0.380** |
| (21%, 58%) | 0.658 | **0.673** | 0.400 | **0.410** | 0.073 | **0.071** | 0.690 | **0.685** |

Table 9: Posterior summaries of parameters for the two spatial models fit to the Restionaceae percent cover data in the Baviaanskloof region. Note that rainfall concentration and minimum annual July temperature were not found to be significant. $\mathcal{M}_2$ includes the spatial effect in $\mu_i$, with the prior for $\delta_0$ centered at 0.5. $\mathcal{M}_3$ includes the effect in $\pi_i$, with the prior for $\gamma_0$ centered at -0.5.

| | $\mathcal{M}_2$ | | | $\mathcal{M}_3$ | | |
|---|---|---|---|---|---|---|
| | mean | 2.5% | 97.5% | mean | 2.5% | 97.5% |
| $\gamma_0$: Intercept | -2.117 | -2.777 | -1.338 | -0.911 | -1.303 | -0.473 |
| $\gamma_1$: mean_ann_pan | - | - | - | 1.167 | 0.054 | 2.367 |
| $\delta_0$: Intercept | 1.029 | 0.762 | 1.465 | 0.543 | 0.201 | 1.045 |
| $\delta_1$: mean_ann_precip | -1.474 | -1.859 | -1.088 | -0.974 | -1.305 | -0.589 |
| $\psi$ | 1.238 | 0.867 | 1.707 | 0.662 | 0.320 | 1.095 |
| $\sigma^2$ | 0.967 | 0.582 | 1.629 | 1.167 | 0.732 | 2.143 |

Table 10: Model comparison metrics for the five models (BEZI, LCEB, $\mathcal{M}_1$, $\mathcal{M}_2$, $\mathcal{M}_3$) fit to Crassulaceae data in Baviaanskloof region. The $\mathcal{M}$ models are fit with $c$ fixed at different values. Best performer for each metric in bold.

| | $R^2$ | AUC | $CRPS_h$ | $CRPS_f$ |
|---|---|---|---|---|
| BEZI | 0.126 | 0.709 | 0.178 | - |
| LCEB ($c = 0.1$) | 0.159 | 0.719 | 0.162 | - |
| LCEB ($c = 0.3$) | 0.174 | 0.718 | 0.123 | - |
| LCEB ($c = 0.5$) | 0.180 | 0.717 | 0.182 | - |
| $\mathcal{M}_1$ ($c = 0.1$) | 0.165 | 0.726 | 0.160 | 0.350 |
| $\mathcal{M}_1$ ($c = 0.3$) | 0.170 | 0.716 | 0.124 | 0.330 |
| $\mathcal{M}_1$ ($c = 0.5$) | 0.172 | 0.720 | 0.089 | 0.299 |
| $\mathcal{M}_2$ ($c = 0.1$) | 0.218 | 0.760 | 0.143 | 0.452 |
| $\mathcal{M}_2$ ($c = 0.3$) | 0.219 | 0.750 | 0.111 | 0.308 |
| $\mathcal{M}_2$ ($c = 0.5$) | **0.225** | **0.766** | **0.080** | **0.277** |
| $\mathcal{M}_3$ ($c = 0.1$) | 0.154 | 0.716 | 0.159 | 0.378 |
| $\mathcal{M}_3$ ($c = 0.3$) | 0.166 | 0.712 | 0.124 | 0.372 |
| $\mathcal{M}_3$ ($c = 0.5$) | 0.156 | 0.714 | 0.089 | 0.358 |

Table 11: Model comparison metrics for the five models (BEZI, LCEB, $\mathcal{M}_1$, $\mathcal{M}_2$, $\mathcal{M}_3$) fit to Restionaceae data in Baviaanskloof region. The $\mathcal{M}$ models are fit with $c$ fixed at different values. Best performer for each metric in bold.

| | $R^2$ | AUC | $CRPS_h$ | $CRPS_f$ |
|---|---|---|---|---|
| BEZI | 0.155 | 0.728 | 0.163 | - |
| LCEB ($c = 0.1$) | 0.341 | 0.833 | 0.192 | - |
| LCEB ($c = 0.3$) | 0.351 | 0.884 | 0.148 | - |
| LCEB ($c = 0.5$) | 0.354 | 0.885 | 0.105 | - |
| $\mathcal{M}_1$ ($c = 0.1$) | 0.289 | 0.867 | 0.193 | 0.404 |
| $\mathcal{M}_1$ ($c = 0.3$) | 0.288 | 0.867 | 0.185 | 0.391 |
| $\mathcal{M}_1$ ($c = 0.5$) | 0.278 | 0.860 | 0.183 | 0.419 |
| $\mathcal{M}_2$ ($c = 0.1$) | 0.367 | **0.908** | 0.148 | 0.415 |
| $\mathcal{M}_2$ ($c = 0.3$) | 0.416 | **0.908** | 0.119 | 0.443 |
| $\mathcal{M}_2$ ($c = 0.5$) | **0.432** | **0.908** | **0.086** | 0.455 |
| $\mathcal{M}_3$ ($c = 0.1$) | 0.278 | 0.871 | 0.161 | **0.303** |
| $\mathcal{M}_3$ ($c = 0.3$) | 0.300 | 0.870 | 0.127 | 0.320 |
| $\mathcal{M}_3$ ($c = 0.5$) | 0.310 | 0.866 | 0.091 | 0.327 |