



Bayesian Analysis of Spatial Zero-Inflated and Right-Censored Survival Data

Sepideh ASADI and Mohsen MOHAMMADZADEH 

The central focus of survival analysis lies in examining the survival time or time to event associated with a preferred outcome. These events are typically characterized by a shift from one distinct state to another. Conventionally, this time is quantified as a continuous, nonnegative real number. However, instances arise where survival data are categorized into discrete time intervals or documented at distinct time points, transforming the survival time into a discrete random variable. When the exact event times are known, employing methods treating time as continuous is appropriate. Conversely, if only the event's day, month, or year is available, discrete-time methods may be more suitable. Discrete-data survival analysis pertains to data values confined to a discrete grid. Notably, discrete-time data sets often exhibit a prevalence of zero values, leading to “zero-inflated discrete survival data.” Despite the increasing utilization of such data across various fields, more research still needs to be done on their spatial modeling and analysis. So, further exploration and refinement of methodologies are required for comprehensive understanding and applicability. This study concentrates on the event of interest, Cercosporiosis disease, seeking to ascertain the annual average time until becoming infected. We aim to scrutinize this data using the zero-inflated discrete Weibull distribution, incorporating a spatially correlated structure. To evaluate the performance of our proposed model, the conditional two-part spatial survival with random effects model will undergo scrutiny with two other models, the classic and spatial model with independent random effects, through simulation studies. The results suggest that the correlated spatial CZIDW model is preferable, especially with a significant correlation between the two parts of the model. Subsequently, we will demonstrate how our model can effectively analyze real-world spatial zero-inflated discrete time-to-Cercosporiosis disease infestation in olive trees. In the garden under our study, the probability of a tree becoming infected with a fungal disease in less than a week is spatially dependent on the average number of weeks until the disease is contracted, provided that it does not fail prematurely. We applied our proposed model to this data set. Evidence showed that using a BICAR prior in a spatial model has enhanced the fit of a spatial model when random effects are correlated instead of considering using an intrinsic CAR model.

Key Words: Survival data; Zero-inflated; Conditional spatial modeling; Random effects; Two-part regression model.

S. Asadi · M. Mohammadzadeh (✉) Department of Statistics, Tarbiat Modares University, Tehran, Iran
(E-mail: mohsen_m@modares.ac.ir).

© 2025 International Biometric Society
Journal of Agricultural, Biological, and Environmental Statistics
<https://doi.org/10.1007/s13253-025-00682-w>

Published online: 07 March 2025

1. INTRODUCTION

Discrete-time survival analysis, introduced by [Cox \(1972\)](#), is a statistical method designed to consider situations where the length of time until a specific and well-defined event, or more generally, time to event occurs, is measured at distinct time points rather than continuously. This means the time until a particular event occurs can only take values from a discrete grid (e.g., 0, 1, 2, ...). It has diverse applications in practical lifetime studies. It can be employed to analyze various scenarios, including the number of treatment sessions until a cure, the number of days before death in certain diseases, the number of equipment runs before failure, the number of hospitalization days, the years remaining until retirement, and the number of students' semesters. When dealing with discrete survival data, there are two primary approaches: One treats the event time as inherently discrete or "truly discrete" ([Jenkins 2005](#)), suitable for scenarios like attempts to solve a puzzle or machine breakdowns ([Nakagawa and Osaki 1975](#)), while the other deals with interval censoring within continuous time frames. Time is divided into distinct chunks in the latter, and discrete values are associated with each interval's beginning, end, or center ([Allison 1982](#)). For instance, values like 0, 1, ... represent the rounding of time intervals with fixed boundaries, such as $[0 - 1)$, $[1 - 2)$, ... The duration of these intervals can be measured in months, years, or even days. The Poisson distribution is commonly employed to model discrete count data due to its straightforward interpretation; however, it has limiting assumptions, such as equal mean and variance, which may not hold in real-world scenarios involving over-dispersion or under-dispersion. This can result in inaccurate inferences and underestimated standard errors. The negative Binomial distribution is often used to address over-dispersion but is unsuitable for under-dispersion. Various generalizations of the Poisson model have been developed for sparsely distributed data, but they still adhere to Poisson distribution assumptions ([Sturman 1999](#)). The Poisson distribution provides the probability of an event occurring a specific number of times within a defined period and assumes the independence of event occurrences. As a result, it does not directly correspond to the time-to-event distribution. Therefore, it becomes essential to generalize discrete distributions in survival analysis that reflect all types of dispersion and do not rely on the assumption of independence, particularly when data are correlated. The proposed type *I* discrete Weibull distribution by [Nakagawa and Osaki \(1975\)](#), designed to mirror the continuous survival function, is well suited for discrete survival data and adept at addressing both over-dispersion and under-dispersion, as highlighted by [Chanialidis et al. \(2017\)](#) and [Kalktawi et al. \(2018\)](#). However, its applicability could be better when confronted with complex time-to-event data characterized by an abundance of zeros. These additional zeros often occur when specific units do not conform to the two basic assumptions of survival analysis. Failing to consider these deviations may lead to two sources producing zero observed values in the same distribution ([Moria et al. 2022](#)). One assumption is that all units experience the event of interest for a sufficient duration. However, some units (immune subjects) may not, leading to structural zeros ([Lawless 2011](#)). Another assumption is that the probability of the event not occurring at time zero is one.

Nevertheless, in particular studies, a subset of subjects may exhibit a notable phenomenon known as an “early event” or fast event. These cases involve a rapid occurrence of the event of interest within a relatively short period after the study commences. Consequently, these subjects may exhibit survival times that are either zero or close to zero, commonly called sampling zeros. The grouping approach of discrete survival data is considered zero when the survival time falls within the first-time category $[0, a_1)$ and is less than the chosen time unit (daily, monthly, yearly, etc.). This situation is common in various fields, such as healthcare, where the number of days in the hospital before delivery for pregnant women is often less than a day. Similarly, when studying divorce or job-seeking periods, a zero value indicates a marital lifetime of less than one year for a couple or immediate success in finding a job, respectively. However, when there is an unexpectedly high number of zero observations, we encounter zero inflation. In the research literature, researchers have designed “cure models” and “zero-inflated survival models” to overcome the limitations of standard models in survival data in accounting for non-susceptible groups and “early event” groups, respectively. For instance, [Braekers and Grouwels \(2016\)](#) used a zero-inflated Cox proportional hazard model to study the sleep time of rats after ethanol consumption. [Oliveira et al. \(2017\)](#) and [Louzada et al. \(2018\)](#) utilized continuous Weibull distribution models to investigate the time to fraud in banking data. [Calsavara et al. \(2019\)](#) extended the defective model to examine time-to-endoscopic stent occlusion and time-to-start insulin use. They proposed a defective regression model incorporating early failures or zero-adjusted outcomes in survival data modeling. [Souza et al. \(2022\)](#) developed a log-normal zero-inflated cure regression model to analyze the labor duration in African pregnant women. Similarly, [Souza et al. \(2022\)](#) discussed a Bayesian estimation approach, specifically the log-normal zero-inflated cure method, to study women diagnosed with invasive cervical cancer in Brazil. [Moria et al. \(2022\)](#) illustrated the use of zero-inflated Bernoulli models, employing a Bayesian approach, to distinguish between different sources of zeros in sickness presenteeism data in the occupational health field. Despite significant advancements in this field, it is crucial to acknowledge that when analyzing zero-inflated discrete time-to-event data, the dispersion structure within the entire dataset may be different from the dispersion patterns observed in the nonzero count group ([Tin 2008](#)). Moreover, the process of right-censoring contributes to the reduction of over-dispersion in the data ([Allison 1982](#)). As soon as large values are censored, the range of observed values becomes narrower, potentially resulting in an underestimation of dispersion, while the data are expected to be overdispersed. Therefore, employing an appropriate zero-inflated distribution can effectively model survival data in such scenarios which becomes imperative. In this paper, we aim to utilize the “Censored ZIDW” (CZIDW) model for analyzing zero-inflated discrete survival data with right-censored observations, which offers the advantage of capturing various dispersion forms in zero and nonzero data.

Incorporating spatial information into the survival model can yield significant benefits, particularly in accounting for variations in survival times across different locations and identifying high-risk areas. [Lawson et al. \(2016\)](#) and [Motarjem et al. \(2020\)](#). The spatial correlation of survival data manifests as the tendency for nearby locations to exhibit similar survival outcomes, a phenomenon commonly observed in diverse survival studies. This correlation arises from shared characteristics among regions, particularly those nearby, and can affect the assumption of independence and data dispersion. To address this spatial cor-

relation structure, specific methodologies, such as conditional (structural) spatial modeling, have been proposed (Banerjee and Dey 2005; Hennerfeind et al. 2006). These approaches are termed “conditional” because the interpretation of regression coefficients is contingent upon spatial frailties (random effects) (Schnell et al. 2015). Hierarchical Bayesian models have been widely utilized to integrate area-level spatial random effects, mainly when data are collected using area (lattice) approaches and precise geographic coordinates (e.g., latitude and longitude) are unavailable. Within this framework, a prior distribution based on the adjacency matrix of areas is employed to integrate information from neighboring regions and account for the correlation between strata concerning their proximity. The extension of these methodologies to zero-inflated survival data is crucial due to the potential for a high prevalence of zeros resulting from a solid spatial correlation in count data, as highlighted by Lee et al. (2016). This phenomenon arises from repeating zero survival times in numerous samples from adjacent regions, attributable to the similarities between the geographical areas. In the realm of zero-inflated area-referenced data, Zhu et al. (2015) has developed Bayesian two-part spatial models, encompassing binary components for the proportion of structural zeros (i.e., the likelihood of observing a zero survival time) and count components for nonzero values in a susceptible population [i.e., discrete time-to-event values (Nyandwi et al. 2020)], to address spatial correlation effectively. In this context, the spatial effect is incorporated either by including a Gaussian random effect via a univariate conditional autoregressive (CAR) distribution prior in the nonzero part (Ver Hoef and Jansen 2007) or by integrating two Gaussian spatial random effects with a bivariate conditional autoregressive distribution (BICAR) in both models (Neelon et al. 2015) to account for spatial correlation. In a BICAR prior, the variance–covariance matrix enables the modeling of any level of dependence between the two components of the model. While many studies have been in the literature on spatial modeling of zero-inflated data for discrete or semi-continuous data, the issue of modeling “censored zero-inflated-discrete time-to-event” data remains unaddressed. Therefore, inspired by Neelon et al. (2015), this paper aims to extend the conditional two-part spatial model to handle zero-inflated discrete survival data with right-censored observations following the CZIDW. We employ a bivariate conditional autoregressive prior distribution to capture the spatial random effects. In our approach, we adopt a Bayesian modeling framework and develop a practical computational algorithm based on the Markov chain Monte Carlo method. This algorithm primarily relies on easily sampled Gibbs steps, ensuring efficiency and effectiveness for our purposes. Finally, we apply our proposed model to analyze the spread dynamics of Cercosporiosis, a fungal disease impacting olive groves in western Iran, focusing on the disease’s early temporal progression following infection. We investigate the unexpected increase in the number of infected trees during the initial weeks by dividing the time frame into nine one-week intervals: $[0 - 7)$, $[7, 14)$, \dots , $[60, \dots)$. The overall outline of the paper is as follows: Sect. 2 introduces the dataset and presents the exploratory analysis that encourages the use of zero-inflated models. Section 3 describes some fundamentals of zero-inflated, censored, discrete Weibull survival analysis concepts. Section 4 introduces spatial models that describe the spatial variation in zero-inflated and right-censored survival data, such as using random effects. Section 5 describes an efficient Bayesian parameter estimation method. Section 6 introduces a simulation study to compare three models: the CZIDW model, the spatial CZIDW model

with independent random effects, and the spatial CZIDW model with correlated random effects. A comparison of the methods for different zero-inflated model for the real dataset is provided in Sect. 7. Section 8 summarizes the results.

2. DESCRIPTION OF THE DATA

Crop cultivation requires careful management of agricultural practices to maximize yields. This involves assessing crops, selecting suitable seed varieties, identifying necessary nutrients, and considering agronomic factors. Crop diseases pose challenges to productivity and food security, with plant diseases disrupting normal functions and affecting productivity. Effective management protects yields from losses due to pathogens and natural disasters. Proactive measures are crucial to identify and address diseases that may impact harvests. Plant diseases can hinder growth rates, reduce fruit production, and spread through various mechanisms. Notably, olive producers in Iran face significant challenges from diseases like olive leaf spot. This paper analyzes the prevalence of Cercosporiosis in an olive grove in western Iran caused by the fungal agent *Cercospora*. Cercosporiosis is a fungal disease characterized by the formation of greenish-brown sclerotial spots that are loosely circular or somewhat cylindrical. This fungus primarily affects olive trees, causing visible spots on the leaves. The symptoms of Cercosporiosis include gray to lead-colored spots on the lower surface of the leaves, which eventually turn completely black. As the disease progresses, the upper surfaces of the leaves may also exhibit yellowing, leading to premature leaf drop. Figure 1 illustrates these symptoms in olive trees. Timely and accurate diagnosis of Cercosporiosis is crucial from the moment the disease first appears. This proactive approach not only enhances disease management strategies but also contributes to the overall sustainability and productivity of olive cultivation. However, there are significant gaps in our understanding of the disease's epidemiology and transmission timeline, particularly regarding the spatial analysis of how the disease spreads over time and space.

To investigate the spread dynamics of Cercosporiosis, this study focuses on a 5000-square-meter olive garden containing 2200 saplings. For each affected tree, the time of infection onset will be meticulously recorded over a two-month period. The study incorporated three explanatory variables: the age of the trees, the species of olive trees (which included three distinct species), and the height of the trees. For any tree exhibiting signs of infection, as defined by the specific symptoms associated with the disease, the time of infection was meticulously recorded. The findings of this study are of significant importance in understanding the epidemiology of this disease. Cercosporiosis, caused by the fungal agent *Pseudocercospora cladosporioides*, can spread between olive trees relatively quickly, often within a few days under favorable conditions. While the exact transmission time can vary based on environmental factors such as humidity and temperature, it is generally accurate to say that there is a reasonable probability that olive trees near each other may contract the disease in less than a week. This highlights the importance of monitoring and managing olive groves to prevent the spread of this disease effectively.

To further investigate the inflation observed in the time to become infected with *Cercospora* olive trees data, where the number of trees infected within a week is higher than



Figure 1. Cercosporiosis symptoms in an olive orchard.

expected, we have partitioned the time axis into nine 1-week periods, represented by intervals such as $[0 - 7)$, $[7, 14)$, \dots , $[60, \dots)$. These intervals are defined by discrete survival time values, namely 0, 1, and so on, serving as the starting points for each interval. Consequently, the data, characterized by excessive zero counts, are called zero-inflated. In this study, we employ survival analysis techniques, a statistical method commonly used in epidemiology and plant pathology to analyze time-to-event data, to analyze our data. These techniques are beneficial when studying the spread of diseases or the survival of organisms over time. Although the precise geographical coordinates for the units are accessible, our objective is to implement our model for areal data, specifically in the context of lattice data structures. To facilitate this transition, we have segmented our geostatistical dataset into twenty units, each containing nine olive trees. This approach allows us to convert our data into a more structured format suitable for our analysis. After the study, 85 of the trees were affected by the disease, while the remaining trees remained healthy, resulting in a right-censoring rate of approximately 51%. The Kaplan–Meier plot (Kaplan and Meier 1958) is drawn in Fig. 2 to provide visual insights into how the number of infected trees relates to the time to event of the disease ($t_1 < \dots < t_9$). The shaded areas around the curves represent the confidence intervals. The risk table provides additional information about the number of trees at risk at different time points, broken down by strata. This can help understand the number of trees observed at each time point and how it changes over time.

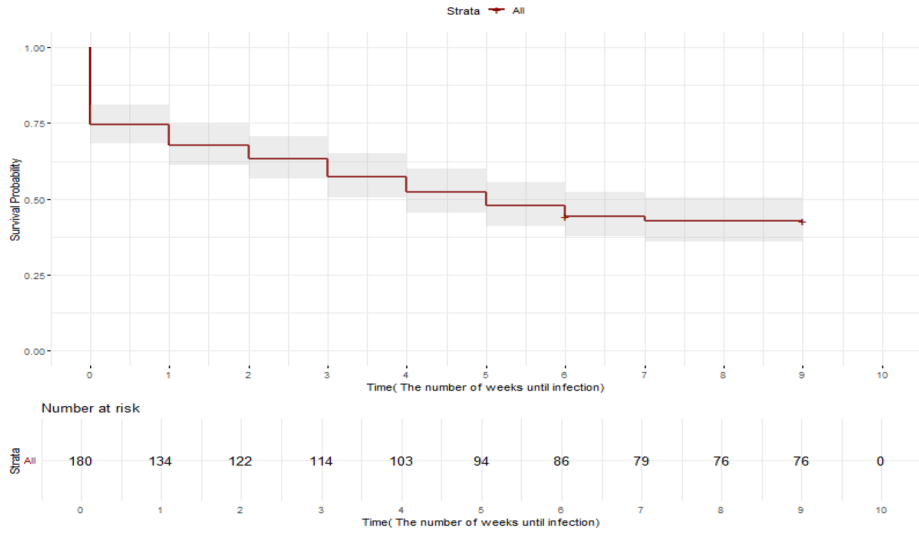


Figure 2. Kaplan–Meier survival curve for “time-to-infestation” in olive trees .

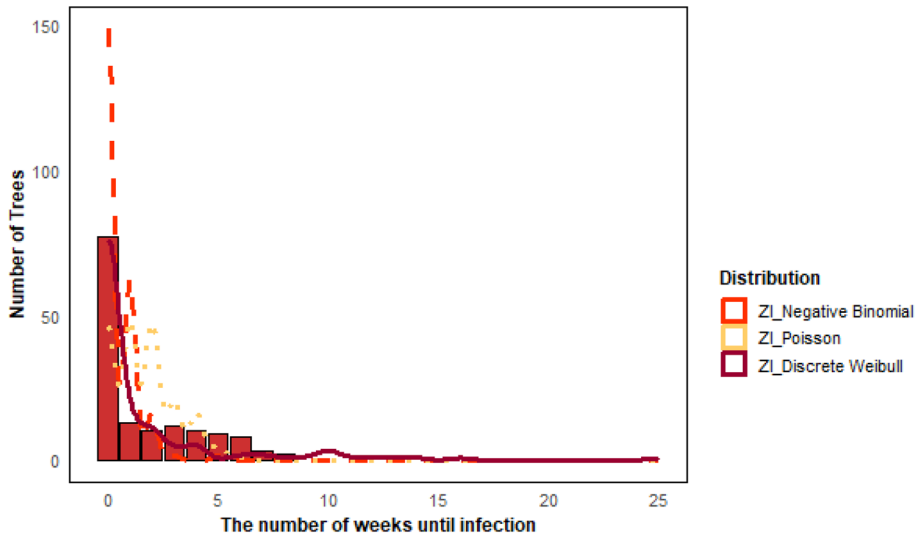


Figure 3. Bar plot with density curve of zero-inflated distributions for “time-to-infestation” in olive trees.

Figure 3 illustrates that around %41 of the infected trees showed signs of *Cercospora* within a week. This suggests the possibility of zero inflation within the data. Consequently, it is imperative to employ censored zero-inflated discrete distributions to model the data. The dispersion index represents the ratio of the observed variance from the data to the observed mean. In this case, the dispersion index equals 2.81, indicating over-dispersion in the data. Three distributions, namely zero-inflated discrete Weibull (ZIDW), zero-inflated negative binomial, and zero-inflated Poisson, have been fitted to the time to *Cercosporiose* disease infestation data. Based on the observation in Fig. 3, it is evident that the ZIDW is better suited for these data than the other two distributions.

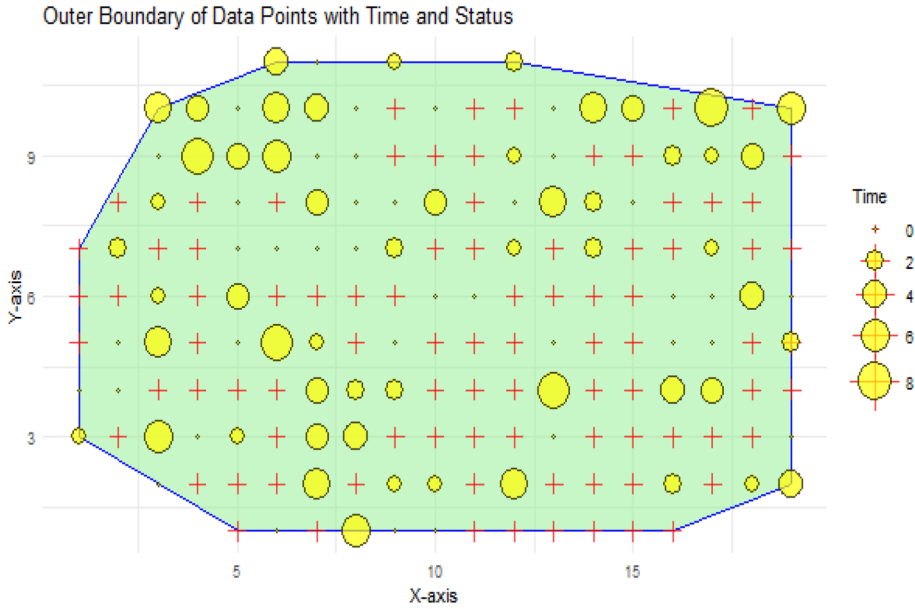


Figure 4. Location and time of infection of trees in the olive orchard, observations (.), and censored data (+).

The mechanisms underlying spore dispersal play a pivotal role in the transmission of diseases within garden ecosystems. Precipitation facilitates the movement of spores from infected leaves to healthy ones within the canopy, while wind carries spores from decaying leaves to nearby trees. These processes underscore the profound impact of geographical factors on the spread of disease. In particular, variables such as tree density and the spatial arrangement of trees within a grove significantly affect the overall dissemination of the disease. Moreover, these factors reveal a clear indication of spatial autocorrelation, suggesting that the risk of infection for any given tree increases as its distance to already infected trees decreases. Figure 4 illustrates the spatial distribution of the trees in the garden; affected trees are denoted with a circle, whereas those that remained healthy by the end of the study are indicated with a plus sign. The size of the yellow circles correlates with the timing of the disease manifestation, with smaller circles representing earlier infections observed within the initial weeks of the study. This visual representation aids in comprehensively understanding the spatial and temporal dynamics of the disease spread within the garden ecosystem.

To effectively analyze the spatial dynamics of our geostatistical dataset and facilitate the examination of lattice data, we aggregated tree locations into 20 squares, each containing nine trees. This aggregation simplified the spatial correlation structure and established transparent spatial relationships between adjacent squares. Given the presence of zero inflation in our dataset, we computed two essential components for each square: the probability of zero survival time (the likelihood of infection occurring within the first week) and the mean nonzero survival times, which represent the average lifespan of trees that were not infected during this initial period. Figure 5 visually illustrates the geographical variation of these two components across the square-segmented areas within the garden. The observed similarities in patterns across neighboring segments strongly indicate significant spatial autocorrelation.



Figure 5. The spatial pattern of a: The probability of Cercorpiose disease infestation in less than a week, b: The average duration of survival time for trees with a lasting lifetime of more than a week .

Specifically, the square segments close to each other on both maps exhibit similar conditions. This observation allows us to implicitly conclude that there is a spatial correlation between the observations for both components, reinforcing the interconnectedness of the spatial dynamics in our dataset.

To quantitatively evaluate spatial autocorrelation, Moran’s statistics (Moran 1950) were calculated to both components, yielding significant values of 0.30 for the “probability of early Cercorpiose disease infestation” (the probability of being zero) component and 0.47 for “the average of non-zero values” component. These results indicate a considerable positive spatial autocorrelation. When comparing the obtained p -values of 0.0155 and 0.0006 against the threshold of significance, both were found to be significant at the ($\alpha = 0.05$) level. This research carries substantial implications, as the null hypothesis of no spatial correlation is firmly rejected, confirming the existence of spatial correlation in both components as a definitive reality, thereby challenging established theories. Additionally, these two components exhibit a notable correlation of -0.42 . This correlation highlights the relationship between the probability of early Cercorpiose disease infestation and the average “time to infestation” across various square segments of the orchard. These spatial patterns, which we will discuss in greater detail in the following section, are essential for understanding the dynamics of Cercosporiosis disease infestation in olive trees.

The Cercosporiosis disease, caused by the fungus *Cercospora*, can be influenced by several factors that affect the timing and severity of infestation; in this study, we conducted an analysis that took into account several covariates, including the age (in years), type (tree distinct varieties), and height (in meters) of each olive tree, to include in the CZIDW model as auxiliary variables according to Table 1.

3. ZIDW DISTRIBUTION

Weibull distribution has a lot of flexibility in modeling survival data with different risk rates (ascending, descending, and fixed). If the random variable T follows the discrete Weibull distribution of the first type (Nakagawa and Osaki 1975), $T \sim DW_I(q, \beta)$, then its density and cumulative distribution functions are, respectively, given by

Table 1. Demographic characteristics

Variable	Group	Number	Percent
Type	I	43	0.23
	II	84	0.46
	III	53	0.29
Variable	Mean	Standard deviation	Range
Age	4.50	0.83	(2.60, 6.10)
Height	3.90	0.84	(2.80, 6.60)

$$\begin{aligned}
 f(t) &= P(T = t) = q^{t^\beta} - q^{(t+1)^\beta}, \quad t = 0, 1, 2, \dots \\
 F(t) &= P(T \leq t) = 1 - q^{(t+1)^\beta}.
 \end{aligned} \tag{1}$$

where $q \in (0, 1)$. The small q means that an excessive zero case occurs. In (1), $\beta > 0$ is another shape parameter that controls the skewness of the distribution. If $\beta \rightarrow \infty$, the discrete Weibull distribution will lead to the Bernoulli distribution, and if $\beta \rightarrow 0$, the distribution will be highly skewed (Kalktawi et al. 2018). The survival and hazard functions corresponding to this distribution are, respectively, given by

$$\begin{aligned}
 S(t) &= P(T \geq t) = 1 - F(t) = q^{(t)^\beta}, \\
 h(t) &= \frac{f(t)}{S(t)} = q^{(t)^\beta - (t+1)^\beta} - 1.
 \end{aligned}$$

Some of the primary discrete Weibull model's characteristics indexes, such as dispersion (DI), zero inflation (ZI), and heavy-tail (HT), are defined about the Poisson distribution as follows

$$DI = \frac{VAR(T)}{E(T)}, \quad ZI = 1 + \frac{\log(P(T = 0))}{E(T)}, \quad HT = \frac{P(T = t)}{P(T = t + 1)}, \quad t \rightarrow +\infty$$

The dispersion index DI denotes over-dispersion, under-dispersion, and equip-dispersion for values of $DI > 1$, $DI < 1$, and $DI = 1$, respectively. The zero inflation ZI index denotes zero inflation for $ZI > 0$, zero deflation for $ZI < 0$, and no excess of zeros for $ZI = 0$. The heavy-tail index HT indicates a distribution for $HT \rightarrow 1$ when $T \rightarrow 1$. These three characteristics for different values of parameters of discrete Weibull distribution show the high flexibility of this distribution.

3.1. ZIDW REGRESSION MODEL

Let the survival time T be a nonnegative random count variable with zero-inflated distribution

$$T \sim \begin{cases} 0 & \text{with probability } \pi \\ f_1(t) & \text{with probability } 1 - \pi \end{cases}$$

In this model, it is presumed that for each observation, there are two data generation processes with different probabilities: One with a chance of $0 < \pi < 1$ creates zero, and the other with a case of $1 - \pi$ makes counting numbers that follow the $f_1(t)$ distribution. A Bernoulli model can determine which of the two processes is used. In this case, $f_T(t) = \pi 1_{(t=0)} + (1 - \pi)f_1(t)$. So we have

$$T = \begin{cases} 0 & \text{with probability } \pi + (1 - \pi)f_1(0) \\ k & \text{with probability } (1 - \pi)f_1(k) \end{cases}$$

which is called a mixed density function with the parameter π , and its corresponding distribution function is given by $F_T(t) = \pi + (1 - \pi)F_1(t)$. The ZIDW distribution can model highly skewed count data with more zeros, thus a good candidate for zero-inflated survival data. If $f_1(t)$ is a discrete Weibull distribution with mean μ , we have

$$\begin{aligned} f_1(0) &= P(T = 0) = \pi + (1 - \pi)(1 - q), \\ f_1(k) &= P(T = k) = (1 - \pi)(q^{k^\beta} - q^{(k+1)^\beta}) \quad k = 1, 2, 3, \dots \end{aligned} \quad (2)$$

The ZIDW model reduces to a standard discrete Weibull distribution if $\pi = 0$. When we want to assess the effect of covariates on the response variable that follows the discrete Weibull distribution, we use the discrete Weibull regression model with linked functions of parameters q and β . In this paper, we connected the covariates only through the parameter q . Regarding the continuous Weibull model, q is equal to $\exp(-\lambda)$ (Khan et al. 1989), where λ is the scale parameter in the Weibull continuous distribution, and then, to create a regression model, it is assumed that the parameter λ is related to covariates through a generalized linear regression model $\log(\lambda) = X'\alpha$, where $X_{n \times (p_1+1)} = (1, X_1, \dots, X_{p_1})$ is a matrix of covariates and $\alpha = (\alpha_0, \dots, \alpha_{p_1})$ is the vector of regression parameters (Da Silva et al. 2008). Consequently, via continuous distribution, the covariates can be arranged with the parameter q of the discrete Weibull distribution as follows: Assume $T|X \sim DW(q(X), \beta)$ that $q(X)$ represents the relation of the parameter q to the covariates, in which case the link function to X would be as follows

$$\log(-\log(q(X))) = X'\alpha, \Rightarrow q \equiv q(X) = e^{-e^{X'\alpha}}. \quad (3)$$

Henceforth, the probability mass function of discrete Weibull distribution is given by

$$f_1(t|x) = \left(e^{-e^{X'\alpha}}\right)^{t^\beta} - \left(e^{-e^{X'\alpha}}\right)^{(t+1)^\beta}, \quad t = 0, 1, \dots$$

If $T|X, Z \sim ZIDW(q(X), \beta, \pi(Z))$, both parameters π and q can depend on covariates that, in the particular case, $\log(q)$ and $\text{logit}(\pi)$ are both linear functions of covariates, and we have

$$\text{logit}(\pi) = \log\left(\frac{\pi}{1 - \pi}\right) = Z'\gamma, \Rightarrow \pi = \frac{e^{Z'\gamma}}{1 + e^{Z'\gamma}} = \left(1 + e^{-Z'\gamma}\right)^{-1}. \quad (4)$$

where $Z_{n \times (p_2+1)} = (1, Z_1, \dots, Z_{p_2})$, $\boldsymbol{\gamma} = (\gamma_0, \dots, \gamma_p)$ is the vector of regression coefficients corresponding to $\text{logit}(\pi)$. Hence,

$$f_1(0) = \frac{1}{1 + e^{Z'\boldsymbol{\gamma}}} (1 + e^{Z'\boldsymbol{\gamma}} - e^{-e^{Z'\boldsymbol{\gamma}}})$$

$$f_1(k) = \frac{1}{1 + e^{Z'\boldsymbol{\gamma}}} \left(e^{-k^\beta e^{X'\boldsymbol{\alpha}}} - e^{-(k+1)^\beta e^{X'\boldsymbol{\alpha}}} \right), \quad k = 1, 2, \dots$$

4. SPATIAL ZIDW MODELS

In survival analysis, spatial referencing is essential in understanding the outcome variation. Although with the growing availability and utilization of geographic information systems software, researchers now have enhanced access to georeferenced locations of subjects in their studies, in time-to-event analysis studies, spatial data routinely collected at area levels and divided into a limited number of regional units with clear boundaries, such as counties or some other administrative delineation. Such data are often referred to as areal data.

In these scenarios, the application of spatial survival analysis becomes particularly valuable as it identifies geographic areas exhibiting high and low prevalence of the event of interest. By employing spatial survival analysis techniques, researchers can effectively uncover spatial patterns and trends, thereby facilitating a more comprehensive understanding of the underlying factors influencing the occurrence and distribution of the event of interest. Due to spatial areas that can borrow knowledge from their neighbors and share similar environmental and social factors, adjacent areas have more impact on each other, which leads to the spatial correlation between them. In this case, area-level spatial models are needed to calculate the potential correlated relationship between neighbor areas. Spatial correlation in survival data causes the independence assumption, which is one of the main assumptions of this model, not to be established in investigating the effect of covariates on survival by popular semi-parametric models in survival studies, such as Cox proportional hazard. To consider spatial correlations, we need valid inference on the association of the covariates with the survival times through spatially random effects models. These methods are generally called conditional forms (structural model) or classified directly (non-structural model). In conditional models, random effects, often referred to as frailty, are integrated within the frameworks of the proportional hazards model (Banerjee et al. 2004; Hennerfeind et al. 2006) and the proportional odds model (Banerjee and Dey 2005). These random effects are employed to address spatial correlation within the model. Then, a Gaussian spatial random field is used to consider the spatially correlated frailty factor. Then, they use a Gaussian spatial random field to justify the spatially correlated frailty factors. These methods are called “conditional” because the regression coefficients are interpreted conditional on spatial frailty (Schnell et al. 2015). Extending this method for spatial zero-inflated survival data with right-censored observations is vital for several reasons. Firstly, right-censoring plays a crucial role in reducing data dispersion. Secondly, a high spatial correlation in the data causes a more significant inflation of observations at zero points. Thirdly, the level of

dispersion differs between the group of zero data and the nonzero counting data. Lastly, spatial correlation can significantly increase the number of right-censored observations.

4.1. CONDITIONAL SPATIAL ZIDW MODEL WITH RANDOM EFFECTS

Assume that A_1, \dots, A_L , are a finite number of separated areas, each with the number of units n_1, \dots, n_L , respectively, and $T_{i\ell}|X, Z \sim ZIDW(\pi_{i\ell}(Z), q_{i\ell}(X), \beta)$ is the survival time for observation i at ℓ th area. We now extend the Models (4) and (3) to accommodate spatially autocorrelation structure in data by introducing spatial random effects for area data. In ZIDW models, it is assumed that there are two processes. Process 1 generates only zeros with probability $\pi_{i\ell}$, whereas process 2 generates counts from a discrete Weibull distribution with probability $1 - \pi_{i\ell}$. So we have two components in zero-inflated data, count nonzero and probability being zero $\pi_{i\ell}$. In this situation, to insert the spatial dependence in the form of a spatial random effect into the ZIDW model, we have three approaches:

1. Assumes that spatial correlations of components are accounted for by the unstructured random intercepts ϕ_ℓ only for nonzero components. Thus, the probability of being zero is spatially unrelated.
2. The first approach is not reasonable for our data, because Fig. 5 shows a spatial autocorrelation in the probability of being zero. Therefore in this approach, it is assumed that spatial correlation between outcomes is accounted for in both components by two spatial random effects $\phi_1 = (\phi_{11}, \dots, \phi_{1L})$ and $\phi_2 = (\phi_{21}, \dots, \phi_{2L})$ which are independent. So we have two distinct spatial processes in the model as follows

$$\log(-\log(q(x_{i\ell}))) = \mathbf{x}'_{i\ell}\alpha_{i\ell} + \phi_{2\ell}, \quad \Rightarrow \mathbf{q} \equiv \mathbf{q}(\mathbf{X}) = \mathbf{e}^{-\mathbf{e}^{\mathbf{X}'\alpha + \phi_2}}. \quad (5)$$

$$\text{logit}(\pi_{i\ell}) = \log\left(\frac{\pi_{i\ell}}{1 - \pi_{i\ell}}\right) = \mathbf{Z}'_{i\ell}\gamma_{i\ell} + \phi_{1\ell}, \quad \Rightarrow \boldsymbol{\pi} \equiv \frac{\mathbf{e}^{\mathbf{Z}'\gamma + \phi_1}}{1 + \mathbf{e}^{\mathbf{Z}'\gamma + \phi_1}} = (\mathbf{1} + \mathbf{e}^{-\mathbf{Z}'\gamma + \phi_1})^{-1}, \quad (6)$$

where $\boldsymbol{\gamma}_{i\ell} = (\gamma_{0i\ell}, \dots, \gamma_{pi\ell})$ and $\boldsymbol{\alpha}_{i\ell} = (\alpha_{0i\ell}, \dots, \alpha_{pi\ell})$ are the model components that can be estimated by fitting two separate GLM regression models. So

$$P(T_{i\ell} = k) = \begin{cases} \frac{1}{1 + e^{z_{i\ell}\gamma_{i\ell} + \phi_{1\ell}}} (1 + e^{z_{i\ell}\gamma_{i\ell} + \phi_{1\ell}} - e^{-e^{x_{i\ell}\alpha_{i\ell} + \phi_{2\ell}}}) & k = 0 \\ \frac{1}{1 + e^{z_{i\ell}\gamma_{i\ell} + \phi_{1\ell}}} [(e^{-e^{x_{i\ell}\alpha_{i\ell} + \phi_{2\ell}}})^{k\beta} - (e^{-e^{x_{i\ell}\alpha_{i\ell} + \phi_{2\ell}}})^{(k+1)\beta}] & k = 1, 2, \dots \end{cases}$$

Usually, two spatially structured random effect vectors ϕ_1 are assumed, and ϕ_2 separately follow the intrinsic conditional autoregressive (ICAR) prior distribution (Besag et al. 1991). The ICAR prior distribution is a form of prior distribution utilized in Bayesian spatial modeling. Its frequent application lies in capturing spatial interdependence within data, especially in scenarios involving areal (lattice) data characterized by irregularly shaped spatial units that do not conform to a regular grid representation. The ICAR prior posits that the value of a given spatial unit is conditionally reliant on the values of its neighboring units in which the spatial correlation is addressed in the model through the adjacency matrix. The precision matrix is defined

to capture the spatial dependence structure. Therefore, the conditional distribution is

$$p(\phi_{r\ell} | \phi_{r\ell'}, \ell' \neq \ell, \tau_\ell^{-1}) = N\left(\frac{\sum_{\ell \sim \ell'} \phi_{r\ell'}}{m_\ell}, \frac{\tau_\ell}{m_\ell}\right) \quad (7)$$

where m_ℓ is the number of neighbors of ℓ area. The spatial random variable $\phi_{1\ell}$ for the region ℓ , which has a set of neighbors $\ell' \neq \ell$, whose cardinality is equal to m_ℓ , is normally distributed with a mean equal to the mean of its neighbors. Its variance decreases with the increase in the number of neighbors. The joint distribution (7) rewrites to the pairwise difference formulation:

$$p(\phi | \tau) \propto \tau^{\frac{n-NC}{2}} \exp \left\{ -\tau^2 \sum_{\ell \sim \ell'} \left(\phi_\ell - \phi_{\ell'} \right)^2 \right\}$$

In the context of defining the spatial proximity matrix for areal subregions, NC represents the number of components in the graph. If NC is 1, it means a fully connected areal graph where each subregion is reachable from every other subregion through a sequence of neighbors. Upon examining the pairwise difference formulation, it is clear that the joint distribution lacks identifiability; adding a constant to all elements of ϕ does not alter the joint distribution. This issue is addressed by imposing the constraint that $\sum_\ell (\phi_\ell) = 0$

3. The third approach is a spatial correlation between the components of two random effect vectors $\phi_1 = (\phi_{11}, \dots, \phi_{1L})$ and $\phi_2 = (\phi_{21}, \dots, \phi_{2L})$. To accommodate this potential correlation with the ICAR prior distribution, we assume that $\phi_\ell = (\phi_{1\ell}, \phi_{2\ell})$ from a bivariate ICAR (BICAR) prior distribution (Carlin and Banerjee 2003; Gelfand and Vounatsou 2003), given by

$$\phi_\ell | \phi_{(-\ell)}, \Sigma \sim N_2\left(\frac{1}{m_\ell} \sum_{\ell \in \partial_\ell} \phi_\ell, \frac{1}{m_\ell} \Sigma\right), \quad (8)$$

where m_ℓ represents the number of neighbors in the ℓ region, ∂_ℓ set of neighbors for area ℓ and

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \rho \sqrt{\Sigma_{11} \Sigma_{22}} \\ \rho \sqrt{\Sigma_{11} \Sigma_{22}} & \Sigma_{22} \end{pmatrix}$$

The prior distribution (8) states that the conditional mean ϕ_ℓ is an average of the spatial effects of neighboring regions, with the covariance matrix Σ based on the number of neighbors of the region ℓ of the information prior distribution combines from neighbors via conditional averaging, thus allowing adjacent block groups to effectively “borrow” information from each other. This information sharing can yield more reliable random effect predictions for areas with small sample sizes. Furthermore, the covariance matrix shows that as the number of neighbors increases, more information is available to predict ϕ_ℓ and hence we have more prior confidence that ϕ_ℓ conditionally is similar to the average of its neighbors. As such, the scaled covariance provides a degree of spatial smoothing.

The off-diagonal element Σ_{12} represents the covariance within the region where the spatial dependence relationship between the two effect components controls the spatial randomness of $\phi_{1\ell}$ and $\phi_{2\ell}$. If $\Sigma_{12} = 0$, the two model components are independent and governed by distinct spatial processes. In short, $\phi_{\ell r} \sim BICAR(\bar{\phi}_{\ell r}, \Sigma)$ for $r = 1, 2$ and $\ell = 1, \dots, L$. Like Neelon et al. (2015) and Loquiha et al. (2018), we first obtain the estimate of Σ_{12} and if there is insufficient evidence to conclude $\Sigma_{12} \neq 0$. We can proceed by fitting a reduced model that assumes independent model components.

4.2. SPATIAL CZIDW MODEL WITH RANDOM EFFECTS

Suppose there is a random sample of size $n = \sum_{\ell}^L n_{\ell}$ of survival data. If $T_{i\ell}$ is the exact survival time of the i th unit in the ℓ th area, that is right-censored by fixed constant $C_{i\ell}$, by defining $Y_{i\ell} = \min(T_{i\ell}, C_{i\ell})$, $\delta_{i\ell} = 1$ if $T_{i\ell} \geq C_{i\ell}$ and $J_{i\ell} = 1$, if $y_{i\ell} = 0$. So we can divide data as follows

$$\begin{cases} J_{i\ell} = 1, \delta_{i\ell} = 0 & Y_{i\ell} \text{ is zero and not right-censored} \\ J_{i\ell} = 0, \delta_{i\ell} = 0 & Y_{i\ell} \text{ is nonzero and not right-censored} \\ J_{i\ell} = 0, \delta_{i\ell} = 1 & Y_{i\ell} \text{ is nonzero and right-censored} \end{cases} \quad (9)$$

In this case, the likelihood of the CZIDW model can be defined as follows

$$\begin{aligned} L_{Area} &= \prod_{\ell=1}^L \prod_{i=1}^{n_{\ell}} f_{i\ell} \\ &= \prod_{\ell=1}^L \prod_{i=1}^{n_{\ell}} [P(Y_{i\ell} = 0 | z_{i\ell})]^{J_{i\ell}} [P(Y_{i\ell} = y_{i\ell} | y_{i\ell} > 0, x_{i\ell})]^{1-J_{i\ell}} \\ &= \prod_{\ell=1}^L \prod_{i=1}^{n_{\ell}} ([P(Y_{i\ell} = 0 | z_{i\ell})]^{1-\delta_{i\ell}})^{J_{i\ell}} ([f(y_{i\ell} | x_{i\ell})]^{1-\delta_{i\ell}} [P(Y_{i\ell} \geq C_{i\ell} | x_{i\ell})]^{\delta_{i\ell}})^{1-J_{i\ell}} \\ &= \prod_{\ell=1}^L \prod_{i=1}^{n_{\ell}} \{[\pi_{i\ell} + (1 - \pi_{i\ell})(1 - q_{i\ell})]^{J_{i\ell}} [(1 - \pi_{i\ell})(q_{i\ell}^{y_{i\ell}} - q_{i\ell}^{(y_{i\ell}+1)})]^{1-J_{i\ell}}\}^{1-\delta_{i\ell}} \\ &\quad \times \{1 - [\pi_{i\ell} + (1 - \pi_{i\ell})(1 - q_{i\ell}^{C_{i\ell}})]\}^{\delta_{i\ell}} \end{aligned} \quad (10)$$

and through replacing (5) and (6) in (10), the likelihood of the spatial CZIDW regression model including random effects can be formed as follows

$$\begin{aligned} L_{Area}(\beta, \alpha, \gamma | n, Y, X, Z) &= \prod_{\ell=1}^L \prod_{i=1}^{n_{\ell}} \left(\left[\left(e^{-z_{i\ell}^T \gamma + \phi_{1\ell}} + 1 \right)^{-1} \right] + \left(1 - \left[\left(e^{-z_{i\ell}^T \gamma + \phi_{1\ell}} + 1 \right)^{-1} \right] \right) \left(1 - (e^{-x_{i\ell}^T \alpha + \phi_{2\ell}}) \right) \right)^{J_{i\ell}(1-\delta_{i\ell})} \\ &\quad \times \left[\left(1 - \left[\left(e^{-z_{i\ell}^T \gamma + \phi_{1\ell}} + 1 \right)^{-1} \right] \right) \left(\left(e^{-x_{i\ell}^T \alpha + \phi_{2\ell}} \right)^{y_{i\ell}} - \left(e^{-x_{i\ell}^T \alpha + \phi_{2\ell}} \right)^{(y_{i\ell}+1)} \right) \right]^{(1-J_{i\ell})(1-\delta_{i\ell})} \\ &\quad \times \left\{ 1 - \left[\left(\left(e^{-z_{i\ell}^T \gamma + \phi_{1\ell}} + 1 \right)^{-1} \right) + \left(1 - \left[\left(e^{-z_{i\ell}^T \gamma + \phi_{1\ell}} + 1 \right)^{-1} \right] \right) \left(1 - (e^{-x_{i\ell}^T \alpha + \phi_{2\ell}})^{C_{i\ell}} \right) \right] \right\}^{\delta_{i\ell}} \end{aligned}$$

5. BAYESIAN INFERENCE

To obtain the estimation of the vector model parameters $\theta = (\alpha, \gamma, \beta)$ with the parameter space Θ , we adopt a hierarchical Bayesian modeling approach, which in that the model parameters are considered as random variables, each of which follows a particular prior distribution $\pi(\theta)$. Assuming that the data $\mathbf{B} = (\mathbf{X}, \mathbf{Y}, \mathbf{Z}, \mathbf{J}, \delta)$ are $(y_{11}, X_{11}, Z_{11}\delta_{11}, J_{11}), \dots, (y_{n_L L}, X_{n_L L}, Z_{n_L L}, \delta_{n_L L}, J_{n_L L})$, where for survival time $y_{i\ell}$ and for full data $\delta_{i\ell} = 1$ and for censored data $\delta_{i\ell} = 0$, $\mathbf{x}_{i\ell} = (x_{1i\ell}, \dots, x_{pi\ell})$. In this case, the posterior probability of CZIDW model, using Bayes rule, is

$$P(\theta|\mathbf{B}) \propto P(\theta)P(\mathbf{B}|\theta).$$

The posterior probability with the model considers the following two-part regression model

$$\begin{aligned} \text{logit}(\pi_{i\ell}) &= \mathbf{z}_{i\ell}^T \boldsymbol{\gamma} + \phi_{1\ell}, \\ \log(-\log(q_{i\ell})) &= \mathbf{x}_i^T \boldsymbol{\alpha} + \phi_{2\ell}, \end{aligned}$$

It can improve the prior information of the analysis and make the resulting estimates of the posterior distribution more flexible (Gilks et al. 1995; Hastings 1970; Gamerman et al. 2006). We consider a bivariate CAR (BICAR) prior distribution for $\boldsymbol{\phi}_\ell = (\phi_{1\ell}, \phi_{2\ell})^T$ according to Mardia (1988) in relation(8). Since $\mathbf{M} - \mathbf{A}$ is singular, the joint prior distribution in (8) is improper distribution. But the posterior Φ is proper distribution. Using the prior (8), the conditional priors for $\boldsymbol{\phi}_1 = (\phi_{11}, \dots, \phi_{1L})'$, and $\boldsymbol{\phi}_2 = (\phi_{21}, \dots, \phi_{2L})'$, are, respectively, given by Neelon et al. (2015)

$$\boldsymbol{\phi}_1 | \boldsymbol{\phi}_2 \sim N_L\left(\rho \sqrt{\frac{\Sigma_{11}}{\Sigma_{22}}} \boldsymbol{\phi}_2, (1 - \rho^2) \Sigma_{11}(\mathbf{M} - \mathbf{A})^+\right), \quad (11)$$

$$\boldsymbol{\phi}_2 | \boldsymbol{\phi}_1 \sim N_L\left(\rho \sqrt{\frac{\Sigma_{22}}{\Sigma_{11}}} \boldsymbol{\phi}_1, (1 - \rho^2) \Sigma_{22}(\mathbf{M} - \mathbf{A})^+\right), \quad (12)$$

where $(\mathbf{M} - \mathbf{A})^+$ is a generalized inverse of the rank-deficient matrix $\mathbf{M} - \mathbf{A}$; $\mathbf{M} = \text{diag}(m_1, \dots, m_n)$, and \mathbf{A} is an $n \times n$ adjacency matrix with $a_{ii} = 0$, $a_{i\ell} = 1$ if spatial units i and ℓ are neighbors, and $a_{i\ell} = 0$ otherwise. This conditional prior specification leads to an efficient Gibbs sampling for spatial random effects. Since there are, unfortunately, no conjugate priors in discrete Weibull regression models (Haselimashhadi et al. 2016), we consider non-informative prior distributions on $\boldsymbol{\alpha}$. Improper prior selection is another way to describe non-informative priors, and improper prior distributions may lead to proper posterior distributions. For each $j = 1, \dots, p$, $\alpha_j \in \Re, \gamma_j \in \Re$ and $\beta \in \Re^+$, respectively, the priors $N(\mu_{\alpha_j}, \sigma_{\alpha_j}^2)$, $N(\mu_{\gamma_j}, \sigma_{\gamma_j}^2)$, and inverse gamma $IG(a, b)$ can be suitable choices. Usually, a classical choice for the prior distribution of the covariance matrix $\boldsymbol{\Sigma}$, which includes the correlation between the random effects ϕ_1 and ϕ_2 , is the conjugate inverse Wishart distribution, $IW(v_0, \mathbf{S}_0)$. So, in summary, the prior distributions are as follows

$$p(\beta) \propto IG(a, b) \quad a = b = 0.5,$$

$$\begin{aligned}
 p(\boldsymbol{\alpha}) &\propto N_p(0, \Sigma_{\boldsymbol{\alpha}}), & \Sigma_{\boldsymbol{\alpha}} &= 100I_p, \\
 p(\boldsymbol{\gamma}) &\propto N_p(0, \Sigma_{\boldsymbol{\gamma}}), & \Sigma_{\boldsymbol{\gamma}} &= 100I_p, \\
 p(\Phi \mid \Sigma) &\propto \text{BICAR}(A, \Sigma), \\
 p(\Sigma) &\propto IW(\nu_0, S_0^{-1}), & \nu_0 &= 3, \quad S_0^{-1} = I_2.
 \end{aligned}$$

Then for $\boldsymbol{\theta} = (\boldsymbol{\alpha}, \boldsymbol{\gamma}, \phi_1, \phi_2, \beta, \Sigma)$, we have:

$$\begin{aligned}
 p(\boldsymbol{\theta} \mid \mathbf{B}) &\propto L_{Area}(\beta, \boldsymbol{\alpha}, \boldsymbol{\gamma} \mid n, Y, X, Z) \frac{b^a \beta^{a-1} e^{-b\beta}}{\Gamma(a)} \times e^{-(\boldsymbol{\alpha}^T \Sigma_{\boldsymbol{\alpha}}^{-1} \boldsymbol{\alpha})} \times e^{-(\boldsymbol{\gamma}^T \Sigma_{\boldsymbol{\gamma}}^{-1} \boldsymbol{\gamma})} \\
 &\times |\Sigma|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2} \Phi^T [(\mathbf{M} - \mathbf{A}) \otimes \Sigma^{-1}] \Phi\right\} \\
 &\times |\Sigma|^{-(\phi_0 + L + 1)/2} \times \exp\left\{-\text{tr}(S_0 \Sigma^{-1})/2\right\}, \tag{13}
 \end{aligned}$$

where $|\Sigma|^{-\frac{1}{2}} \exp\{-\frac{1}{2} \Phi^T [(\mathbf{M} - \mathbf{A}) \otimes \Sigma^{-1}] \Phi\}$ is a BICAR joint prior for $\Phi = (\phi_1, \phi_2)$. Since conjugate priors are not available, posterior distributions are often derived from complex models with a large number of parameters. Therefore, it is impossible to specify the Bayes estimates of the model parameters. Also, the integral of (13) does not have a closed form, so the following full conditional distributions can be used to apply Markov chain Monte Carlo (MCMC) algorithms (Hoff 2009; Metropolis 1953).

$$\begin{aligned}
 p(\beta \mid \dots) &\propto \prod_{\ell=1}^L \prod_{i=1}^{n_{\ell}} \left[\left(\left(e^{-e^{x_{i\ell}^T \boldsymbol{\alpha} + \phi_{2\ell}}} \right)^{y_{i\ell}^{\beta}} - \left(e^{-e^{x_{i\ell}^T \boldsymbol{\alpha} + \phi_{2\ell}}} \right)^{(y_{i\ell} + 1)^{\beta}} \right)^{(1 - J_{i\ell})(1 - \delta_{i\ell})} \right. \\
 &\times \left. \left(\left(e^{-e^{x_{i\ell}^T \boldsymbol{\alpha} + \phi_{2\ell}}} \right)^{C^{\beta}} \right)^{\delta_{i\ell}} \times \frac{b^a \beta^{a-1} e^{-b\beta}}{\Gamma(a)} \right]. \tag{14}
 \end{aligned}$$

$$\begin{aligned}
 p(\boldsymbol{\alpha} \mid \dots) &\propto \prod_{\ell=1}^L \prod_{i=1}^{n_{\ell}} \left[1 - \left(e^{-e^{x_{i\ell}^T \boldsymbol{\alpha} + \phi_{2\ell}}} \right) + \left(e^{-z_{i\ell}^T \boldsymbol{\gamma} + \phi_{1\ell}} + 1 \right)^{-1} \times \left(e^{-e^{x_{i\ell}^T \boldsymbol{\alpha} + \phi_{2\ell}}} \right) \right]^{J_{i\ell}(1 - \delta_{i\ell})} \\
 &\times \left[\left(\left(e^{-e^{x_{i\ell}^T \boldsymbol{\alpha} + \phi_{2\ell}}} \right)^{y_{i\ell}^{\beta}} - \left(e^{-e^{x_{i\ell}^T \boldsymbol{\alpha} + \phi_{2\ell}}} \right)^{(y_{i\ell} + 1)^{\beta}} \right)^{(1 - J_{i\ell})(1 - \delta_{i\ell})} \right] \\
 &\times \left(\left(e^{-e^{x_{i\ell}^T \boldsymbol{\alpha} + \phi_{2\ell}}} \right)^{C^{\beta}} \right)^{\delta_{i\ell}} e^{-\frac{1}{2} \boldsymbol{\alpha}^T \Sigma_{\boldsymbol{\alpha}}^{-1} \boldsymbol{\alpha}}. \tag{15}
 \end{aligned}$$

$$\begin{aligned}
 p(\boldsymbol{\gamma} \mid \dots) &\propto \prod_{\ell=1}^L \prod_{i=1}^{n_{\ell}} \left(1 - \left(e^{-e^{x_{i\ell}^T \boldsymbol{\alpha} + \phi_{2\ell}}} \right) + \left(e^{-z_{i\ell}^T \boldsymbol{\gamma} + \phi_{1\ell}} + 1 \right)^{-1} \times \left(e^{-e^{x_{i\ell}^T \boldsymbol{\alpha} + \phi_{2\ell}}} \right) \right)^{J_{i\ell}(1 - \delta_{i\ell})} \\
 &\times \left[\left(1 - \left[\left(e^{-z_{i\ell}^T \boldsymbol{\gamma} + \phi_{1\ell}} + 1 \right)^{-1} \right] \right)^{1 - J_{i\ell}} \right] \\
 &\times \left[\left(1 - \left[\left(e^{-z_{i\ell}^T \boldsymbol{\gamma} + \phi_{1\ell}} + 1 \right)^{-1} \right] \right)^{\delta_{i\ell} J_{i\ell}} e^{-\frac{1}{2} \boldsymbol{\gamma}^T \Sigma_{\boldsymbol{\gamma}}^{-1} \boldsymbol{\gamma}} \right]. \tag{16}
 \end{aligned}$$

It is necessary to consider that full conditional distribution for ϕ_1 depends on the likelihood contribution of all N observations, while for ϕ_2 it relies only on contributions of $N_1 < N$ positive and nonzero observations. This imbalance of the sample size makes the full conditional distribution of the ϕ_1 and ϕ_2 parameters based on the joint prior (8) not very efficient, and it is better to use the conditional priors (11) and (12) separately to obtain full conditional distributions. Therefore, the full conditional distribution of ϕ_1 , ϕ_2 , and Σ is

as follows

$$\begin{aligned}
 p(\phi_1 | \dots) &\propto \prod_{\ell=1}^L \left(\left(1 - \left(e^{-e^{\mathbf{x}^T \boldsymbol{\alpha} + \phi_2}} \right) \right) + \left[\left(e^{-z^T \boldsymbol{\gamma} + \phi_1} + 1 \right)^{-1} \right] \left[e^{-e^{\mathbf{x}^T \boldsymbol{\alpha} + \phi_2}} \right]^{J_\ell(1-\delta_\ell)} \right. \\
 &\quad \times \left\{ 1 - \left[\left(e^{-z^T \boldsymbol{\gamma} + \phi_1} + 1 \right)^{-1} \right]^{1-J_\ell} \right\}^{1-J_\ell} \left\{ 1 - \left[\left(e^{-z^T \boldsymbol{\gamma} + \phi_1} + 1 \right)^{-1} \right]^{\delta_\ell J_\ell} \right\}^{\delta_\ell J_\ell} \\
 &\quad \times e^{-\frac{1}{2} \left(\phi_1 - \rho \sqrt{\frac{\Sigma_{11}}{\Sigma_{22}}} \phi_2 \right)^T \left[(1-\rho^2) \Sigma_{11} (\mathbf{M} - \mathbf{A})^+ \right]^{-1} \left(\phi_1 - \rho \sqrt{\frac{\Sigma_{11}}{\Sigma_{22}}} \phi_2 \right)}
 \end{aligned} \tag{17}$$

$$\begin{aligned}
 p(\phi_2 | \dots) &\propto \prod_{\ell=1}^L \left(1 - \left(e^{-e^{\mathbf{x}^T \boldsymbol{\alpha} + \phi_2}} \right) \right) + \left[\left(e^{-z^T \boldsymbol{\gamma} + \phi_1} + 1 \right)^{-1} \right] \left[e^{-e^{\mathbf{x}^T \boldsymbol{\alpha} + \phi_2}} \right]^{J_\ell(1-\delta_\ell)} \\
 &\quad \times \left[\left(\left(e^{-e^{\mathbf{x}'^T \boldsymbol{\alpha} + \phi_2}} \right)^{y^\beta} - \left(e^{-e^{\mathbf{x}'^T \boldsymbol{\alpha} + \phi_2}} \right)^{(y+1)^\beta} \right) \right]^{(1-J_\ell)(1-\delta_\ell)} \left[e^{-e^{\mathbf{x}'^T \boldsymbol{\alpha} + \phi_2}} \right]^{C^\beta \delta_\ell} \\
 &\quad \times e^{-\frac{1}{2} \left(\phi_2 - \rho \sqrt{\frac{\Sigma_{22}}{\Sigma_{11}}} \phi_1 \right)^T \left[(1-\rho^2) \Sigma_{22} (\mathbf{M} - \mathbf{A})^+ \right]^{-1} \left(\phi_2 - \rho \sqrt{\frac{\Sigma_{22}}{\Sigma_{11}}} \phi_1 \right)}
 \end{aligned} \tag{18}$$

$$\begin{aligned}
 p(\boldsymbol{\Sigma} | \dots) &\propto \prod_{\ell=1}^L | \boldsymbol{\Sigma} |^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \boldsymbol{\Phi}^T \left[(\mathbf{M} - \mathbf{A}) \otimes \boldsymbol{\Sigma}^{-1} \right] \boldsymbol{\Phi} \right\} \\
 &\quad \times | \boldsymbol{\Sigma} |^{-(\phi_0 + L + 1)/2} \exp \left\{ -\text{tr}(\mathbf{S}_\theta \boldsymbol{\Sigma}^{-1})/2 \right\}
 \end{aligned}$$

With a bit of computation in the full conditional distribution, the posterior distribution will be the inverse Wishart as follows

$$\boldsymbol{\Sigma} | \boldsymbol{\Phi} \sim \text{IW}(\nu_0 + L - 1, \mathbf{S}_0 + \mathbf{S}_\Phi), \tag{19}$$

where $\mathbf{S}_\Phi = \boldsymbol{\Phi}^{*T} (\mathbf{M} - \mathbf{A}) \boldsymbol{\Phi}^*$, and $\boldsymbol{\Phi}^*$ is the $L \times 2$ matrix centered at their mean values spatial random effects with two columns $\boldsymbol{\phi}_1 = (\phi_{11}, \dots, \phi_{1L})^T$ and $\boldsymbol{\phi}_2 = (\phi_{21}, \dots, \phi_{2L})^T$.

Posterior Computation Algorithm:

The random walk Metropolis–Hastings algorithm is implemented to derive the model parameters. Let $\boldsymbol{\theta}$ be the current state value and $\boldsymbol{\theta}^*$ be the proposed value generated from the candidate distribution $q(\boldsymbol{\theta}^* | \boldsymbol{\theta})$. The candidate value $\boldsymbol{\theta}^*$ is accepted with probability $p = \min(1, R_\theta)$. Since in random walk Metropolis–Hastings, the candidate distribution is symmetric and depends only on the distance between the current state value and the proposed value, the acceptance rate is equal to

$$R_\theta = \frac{L(\boldsymbol{\theta}^* | \mathbf{B}) \pi(\boldsymbol{\theta}^*)}{L(\boldsymbol{\theta} | \mathbf{B}) \pi(\boldsymbol{\theta})}$$

The steps of this algorithm can be described as follows:

Step 1: Initializing parameters $\boldsymbol{\theta}^{(0)} = (\boldsymbol{\alpha}^{(0)}, \boldsymbol{\gamma}^{(0)}, \phi_1^{(0)}, \phi_2^{(0)}, \beta^{(0)}, \boldsymbol{\Sigma}^{(0)})$

Step 2: Updating:

2-1: For each i and ℓ , we generate $y_{i\ell}$ from the CZIDW distribution with real values.

2-2: In each iteration, we update the Σ matrix from the posterior distribution (19) through Gibbs sampling and the obtained value in the place first. That is, $\theta^{(0)} = (\alpha^{(0)}, \gamma^{(0)}, \phi_1^{(0)}, \phi_2^{(0)}, \beta^{(0)}, \Sigma^{(1)})$.

2-3: To update the vector of random effects $\phi_{1\ell}$ for each $\ell = 1, \dots, L$, assuming that the vector ϕ_2 is known first from the candidate density from the symmetric univariate t distribution centered on the previous value $\phi_{1\ell}$ we get the candidate values $\phi_{1\ell}^*$. Considering the full conditional distribution (17) as $p(B_{i\ell} | \phi_{1\ell}^{(s)}, \theta)$ in the following acceptance ratio, we update the value of $\phi_1 | \phi_2$.

$$\rho_{\phi_{1\ell}} = \frac{p(\phi_{1\ell}^* | \mathbf{B}, \theta)}{p(\phi_{1\ell}^{(s)} | \mathbf{B}, \theta)} = \frac{\prod_{\ell=1}^n \prod_{j=1}^J p(B_{i\ell} | \phi_{1\ell}^*, \theta)}{\prod_{i=1}^n \prod_{j=1}^J p(B_{i\ell} | \phi_{1\ell}^{(s)}, \theta)} \times \frac{\pi(\phi_{1\ell}^*)}{\pi(\phi_{1\ell}^{(s)})}$$

where $\phi_{1\ell}^*$ and $\phi_{1\ell}^{(s)}$ are the candidate and current values, respectively, and $\phi_{1\ell}$ and $\pi(\phi_{1\ell})$ are based on the conditional bivariate CAR prior distribution of $\phi_\ell | \phi_{(-\ell)}$, given in (8).

2-4: To update the vector of random effects $\phi_{2\ell}$ for each $(\ell = 1, \dots, L)$ assuming that the vector ϕ_1 is known, as in the previous step through a Metropolis–Hastings random walk algorithm first selects the candidate values of $\phi_{2\ell}^*$ from the candidate density from the symmetric univariate distribution t centered on the prior value $\phi_{2\ell}$ and considering the full conditional distribution (18) with the new initial value obtained in the previous steps, i.e., $\theta^{(0)} = (\alpha^{(0)}, \gamma^{(0)}, \phi_1^{(1)}, \phi_2^{(0)}, \beta^{(0)}, \Sigma^{(1)})$ as $p(B_{i\ell} | \phi_{2\ell}^{(s)}, \theta)$ in the acceptance ratio, we update the value of $\phi_2 | \phi_1$.

2-5: To ensure identifiability, we apply the zero-sum constraint to ϕ_1 and ϕ_2 .

2-6: To update β by using full conditional distribution (14), assuming that other parameters are known and the normal distribution centered on the current values of β , the acceptance rate of this parameter is calculated. Then, we decide to accept or reject the new sample. Finally, the obtained value in the initial value $\theta^{(0)} = (\alpha^{(0)}, \gamma^{(0)}, \phi_1^{(1)}, \phi_2^{(1)}, \beta^{(0)}, \Sigma^{(1)})$ obtained in the previous steps, we insert, and $\theta^{(0)} = (\alpha^{(0)}, \gamma^{(0)}, \phi_1^{(1)}, \phi_2^{(1)}, \beta^{(1)}, \Sigma^{(1)})$.

2-7: We update the coefficients α by comparing the value of the function (15) in the new initial deal with the value generated from the candidate distribution (t distribution centered on the current value of α) and the accepted value is replaced with the initial one. That is, $\theta^{(0)} = (\alpha^{(1)}, \gamma^{(0)}, \phi_1^{(1)}, \phi_2^{(1)}, \beta^{(1)}, \Sigma^{(1)})$

2-8: Similar to α , we update the γ coefficients by comparing the value of the function (16) in the new initial value and the generated value from the candidate distribution; the new values are $\theta^{(1)} = (\alpha^{(1)}, \gamma^{(1)}, \phi_1^{(1)}, \phi_2^{(1)}, \beta^{(1)}, \Sigma^{(1)})$.

Step 3: For $r = 2, \dots, R$, we iterate the following steps:

3-1: Draw the random error vector ϵ from a multivariate normal distribution with zero mean vector and covariance matrix as a diagonal matrix, where the main diagonal elements are the inverse of the Fisher information matrix. $\epsilon \sim \mathcal{N}(\mu = \mathbf{0}, \Sigma = \text{diag}(I^{-1}(\theta)))$. Then $\theta^* = \theta^{(r-1)} + \epsilon$.

3-2: Calculate $p = \min(1, R_\theta)$.

3-3: Draw u from the uniform distribution, $u \sim U(0, 1)$. If $u \leq p$, accept θ^* and set $\theta^{(r)} = \theta^*$. If $u > p$, reject θ^* and set $\theta^{(r)} = \theta^{(r-1)}$.

Step 4: Remove the first N elements from the generated chain for a burn-in period, with N being chosen to be large enough that the chain reaches its stationary distribution.

Step 5: Thin the chain by recording only every k th value, to reduce the correlation between generated samples.

Step 6: Bayes estimates for the parameters are the average of the final generated sample values.

In the Bayesian approach, model selection is one of the fundamental ways to determine which model fits the data better. We use the deviance information criterion (DIC) (Spiegelhalter 2002) and Watanabe–Akaike information criterion (WAIC) (Watanabe 2010) to evaluate and compare the CZIDW models and the spatial CZIDW models where the lower values of them, the better the model. DIC is defined as $DIC = \hat{p}(D) + \bar{D}(\theta)$, where $\bar{D}(\theta) = E[D(\theta | y)]$ is the posterior mean deviance of the model and is calculated as a measure of how well the model fits the data (goodness of fit) and $\hat{p}(D) = \bar{D}(\theta) - D(\theta)$ is a measure of the complexity of the model (adequate number of parameters), and equal to the difference between mean deviance and deviance at means. The intuition behind the deviance statistic is to examine the improvement in fit produced by the estimation of the parameter vector θ (Darmofal 2006). The Watanabe–Akaike information criterion (WAIC) is a Bayesian model selection criterion that evaluates model fit while accounting for model complexity, making it especially useful for hierarchical models. Key features of WAIC include its foundation in Bayesian principles, which utilize the posterior distribution of parameters, and its inclusion of a penalty for model complexity based on the adequate number of parameters to mitigate overfitting. The calculation of WAIC is expressed as $WAIC = -2 \times (LPPD - pWAIC)$, where LPPD refers to the log pointwise predictive density, and pWAIC denotes the adequate number of parameters, with lower WAIC values indicating better model performance. WAIC and DIC are essential for model selection in statistical modeling, particularly within Bayesian contexts, as they facilitate the choice of models that effectively fit the data while ensuring good generalization of new data.

6. SIMULATION STUDY

To evaluate the validity of the proposed approach for fitting spatial right-censored zero-inflated survival data, we conducted a simulation study for the following three models that become more complicated in order.

Model 1: CZIDW model,

Model 2: Spatial CZIDW model with independent random effects (with CAR prior),

Model 3: Spatial CZIDW model with correlated random effects (with BICAR prior).

1. For the first model simulation, as the basis for other models, we set the number of areas equal to the number of provinces in Iran, $L = 31$. We then generated a sample of size $n_\ell = 50$ for each area $\ell = 1, \dots, L$, so the total number of observations was $n = \sum_{\ell=1}^L n_\ell$. In this study, we assumed X is equivalent to Z and a similar set of covariates affect q and π . Two continuous explanatory covariates were created

from the distributions $N(0, 1)$ and $U(0, 1.5)$. Consequently, $X = (1, X_1, \dots, X_2)$ corresponded to two explanatory covariates, and the intercept was included in the $n \times 3$ design matrix. The temporal component pdf for the i^{th} , $i = 1, \dots, n_\ell$ observation in province ℓ , follows the distribution $T_{i\ell} | X_{i\ell} \sim ZIDW(\pi(X_{i\ell}), q(X_{i\ell}), \beta)$, so two related responses were controlled under two generalized linear models $\text{logit}(\pi)$ and $\log(-\log(q))$. We set initial values for model parameters $\alpha_{\text{real}} = (-2, 0.5, 0.3)$, $\beta = 1.2$, and $\gamma_{\text{real}} = (1, 1.5, -0.2)$. Then, to investigate the right-censoring, we selected 93% of the generated data as the proportion of right-censored data. Therefore, we considered the quantile 93% of data as the censored point $C_{i\ell}$ and as a threshold to cut the simulated sample, such that all values $y_{i\ell} \geq C_{i\ell}$ were re-valued to be equal to $C_{i\ell}$. Also, if the generated time to event $T_{i\ell}$ is not greater than the generated censored time $C_{i\ell}$, we set $\delta_{i\ell} = 1$; otherwise, its value is considered to be zero. To add a zero-inflated feature for each response, first, a random vector from a uniform distribution $U = (u_1, \dots, u_n) \sim U(0, 1)$ is generated, if $u_{i\ell} \leq \pi_{i\ell}$, set $J_{i\ell=0}$ and $Y_{i\ell} = 0$; otherwise, we considered $J_{i\ell=1}$ and generated $Y_{i\ell}$ from discrete Weibull distribution. We generated the outcome data, under the following two generalized linear models

$$\begin{aligned}\text{logit}(\pi_{i\ell}) &= \gamma_{0\ell} + \sum_{p=1}^2 \gamma_{p\ell} x_{p\ell}, \\ \log(-\log(q_{i\ell})) &= \alpha_{0\ell} + \sum_{p=1}^2 \alpha_{p\ell} x_{p\ell}.\end{aligned}$$

We used normal prior distributions $N(0, \sigma_{\alpha_1}^2)$ and $N(0, \sigma_{\alpha_2}^2)$, for regression coefficients α_1 and α_2 , with precision parameters, $\sigma_{\alpha_p}^{-2} \sim T(10^{-5}, 10^{-5})$, $p = 1, 2$. Similarly, for γ_1 and γ_2 the normal priors $N(0, \sigma_{\gamma_1}^2)$ and $N(0, \sigma_{\gamma_2}^2)$, are considered, respectively, with $\sigma_{\gamma_p}^{-2} \sim T(10^{-5}, 10^{-5})$, $p = 1, 2$. Improper uniform priors were assigned for intercepts $\gamma_{0\ell}$ and $\alpha_{0\ell}$. The inverse gamma prior $IG(0.5, 0.5)$ was considered for the shape parameter β .

2. To simulate the second model (spatial CZIDW model with independent effects), we utilized the adjacency matrix, where the number of rows and columns equals the number of provinces in Iran. We determined the Euclidean distance of provinces based on their geographic position. Next, using the ‘‘Distance scheme,’’ in which the neighbors are defined as a function of the Euclidean distance $d_{\ell\ell'}$ of two districts ℓ and ℓ' , each element of the adjacency matrix $W = (w_{\ell\ell'})$ was defined by

$$w_{\ell\ell'} = \begin{cases} \frac{1}{d_{\ell\ell'}^2}, & d_{\ell\ell'} \leq 500 \text{ km}^2, \\ 0, & d_{\ell\ell'} > 500 \text{ km}^2. \end{cases}$$

To illustrate their spatial autocorrelation, we showed the provinces’ autocorrelation graphs in Fig. 6. There is a stronger association between areas situated closer to one another.

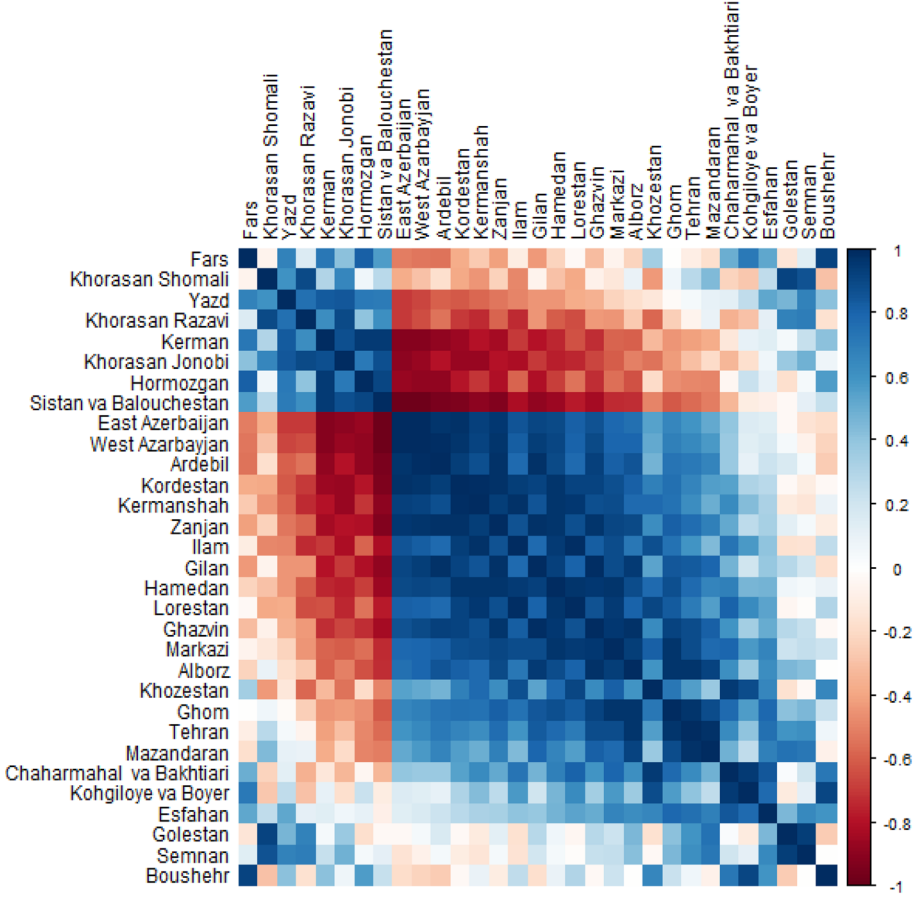


Figure 6. Correlations between provinces of Iran.

Two separate county-level spatial random effects $\phi_{1\ell}$ and $\phi_{2\ell}$ are added in as follows

$$\text{logit}(\pi_{i\ell}) = \gamma_{0\ell} + \sum_{p=1}^2 \gamma_{p\ell} x_{p\ell} + \phi_{1\ell},$$

$$\log(-\log(q_{i\ell})) = \alpha_{0\ell} + \sum_{p=1}^2 \alpha_{p\ell} x_{p\ell} + \phi_{2\ell}.$$

We assigned two independent univariate ICAR priors for $\phi_{1\ell}$, $\phi_{2\ell}$ according to 7, and the uniform priors $U(0, 10)$ and $U(0, 20)$ for the standard deviation of these conditional distributions, respectively. The random effects $\phi_{1\ell}$ and $\phi_{2\ell}$ are latent variables that account for unmeasured area-level factors and, respectively, contribute to the variability in the probability of being zero and nonzero counts of units in the area ℓ . Large values of $\phi_{1\ell}$ indicate a higher occurrence of the event of interest at the zero point in the area ℓ relative to the other areas. In contrast, larger values of $\phi_{2\ell}$ imply a longer time to reach the final event for the ℓ' th region than others.

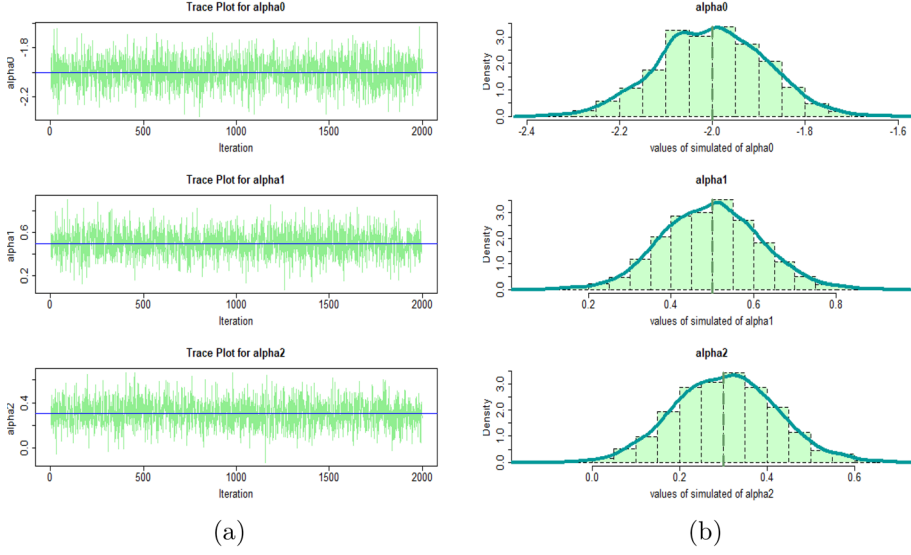


Figure 7. **a** Trace plots, **b** Marginal posterior density plots for the regression coefficients $\alpha_0, \alpha_1, \alpha_2$. The vertical and horizontal dashed lines represent the true values.

- Since regions with a high zero probability may show fewer nonzero counts, we can allow a spatial relationship between two spatial components. Therefore, to simulate the third model (spatial zero-inflated model with correlated random effects), we assigned a BICAR prior for each $\phi_\ell = (\phi_{1\ell}, \phi_{2\ell})^T, \ell = 1, \dots, L$, and the inverse Wishart $IW(v_0, S_0)$ with $v_0 = 3$ and $S_0 = 2I_2$ as a conjugate prior for the covariance matrix Σ . Since the matrix $(\mathbf{M} - \mathbf{A})$ is singular, the spatial random effects cannot be simulated directly. We take $S = 1.38761e - 16$ (Neelon et al. 2015) to avoid this limitation. Now, the spatial random effects $\{\phi_1, \dots, \phi_L\}$ can be generated from the joint prior with $\Sigma = \begin{bmatrix} 4 & 6 \\ 6 & 16 \end{bmatrix}$, and $\rho = 0.75$.

To make objective inferences, we implemented the models using just a single chain to simplify computations. The data-generating processes are repeated 10,000 times with 2000 iterations as burn-in. We retained every 20th observation to reduce autocorrelation. The Gibbs algorithm is applied for each simulated dataset. The MCMC algorithm's convergence was first assessed informally using visual examination by the autocorrelation plots of the traces. Posterior marginal densities and trace plots for α are given in Fig. 7. Kernel density estimates of the marginal posterior distributions are shown on the right, while trace plots are displayed on the left. It shows the convergence and efficiency of the samples produced with the MCMC algorithm to the target distribution and the speed of data correlation reduction. The trace plots show that the MCMC mixes very well. The effect of the starting values for parameters disappears quickly. The shaded area in the density plots represents the 95% credible interval corresponding to the posterior distribution. The true parameter value lies well within its credible interval in all cases. The results indicate the excellent performance of the Bayesian method under ideal conditions.

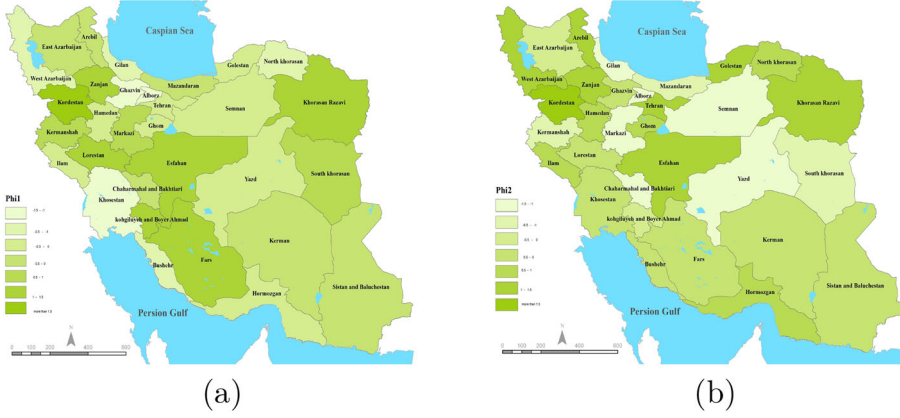


Figure 8. Geographical distribution of Bayesian estimation of spatial random effects a: ϕ_1 , b: ϕ_2 .

The first and second columns of Table 2 display the true values of the regression parameters and the Bayes estimate of parameters of the three model. We compare the methods using the bias and mean squared error (MSE) for each coefficient, calculated as the sum of the square difference between the estimated coefficients at each iteration. We also display the 95% credible interval of the estimated coefficients over the simulations. Table 2 also illustrates how the regression coefficients of all three models have appropriate accuracy. However, in the spatial model with correlated random effects, the estimates of coefficients are closer to the true value. The values of MSE for this model are less than those of the two other models.

Figure 8 displays the geographic distribution of estimated spatial random effects ϕ_1 and ϕ_2 , highlighting significant regional differences in these effects. Dark green provinces are associated with higher estimated values ($\phi_r = 0, r = 1, 2$), while light green areas have lower values ($\phi_r < 0, r = 1, 2$).

As Table 2 indicates, the two spatial models vastly outperformed the other models in terms of fit. While all the coefficients in all models were close to the truth values, most of the coefficient estimates were closer to the truth value in the spatial models than the estimates in the classical model. Correlated random effects were worthwhile, and we evaluate and analyze the link between estimation accuracies and biases with the correlation between $\phi_{1\ell}$ and $\phi_{2\ell}$ four correlation coefficients of 0.75, 0.5, 0.25, and 0, corresponding to the four covariance matrixes $\begin{bmatrix} 4 & 6 \\ 6 & 16 \end{bmatrix}$, $\begin{bmatrix} 1 & 0.5 \\ 0.5 & 4 \end{bmatrix}$, $\begin{bmatrix} 4 & 4 \\ 4 & 16 \end{bmatrix}$ and $\begin{bmatrix} 4 & 0 \\ 0 & 16 \end{bmatrix}$ are evaluated.

Table 3 shows that the two models' parameter estimates, biases, and MSEs were generally similar. However, the bias of the separate model grows as ρ grows. While the random effects are correlated, the independence assumption leads to a positive bias. In addition, for model comparison, we adopt the DIC and p_D , a penalty for model complexity. The penalty p_D was designed to estimate the number of effective parameters in Bayesian hierarchical models. For $\rho = 0$, the DIC value in the spatial model with independent random effects is lower than the same model with correlated random effects, so it is considered preferable. However,

Table 2. Bayesian estimates and 95% credible intervals of the models parameters

Model	Parameter	True value	Estimates	Bias	MSE	%95 credible intervals
1	α_0	-2	-1.989	-0.105	3.857	(-2.016, -1.960)
	α_1	0.5	0.487	-0.012	2.075	(0.458, 0.516)
	α_2	0.3	0.286	-0.013	1.881	(0.255, 0.318)
	γ_0	1	0.995	-0.004	0.064	(0.964, 1.027)
	γ_1	1.5	1.489	-0.010	1.096	(1.159, 1.520)
	γ_2	-0.2	0.207	-0.007	1.490	(-0.238, -0.176)
	Σ_{11}	4	3.992	-0.007	48.897	(3.971, 4.013)
2	α_0	-2	-1.94	0.005	3.868	(-2.026, -1.962)
	α_1	0.5	0.487	-0.0127	2.136	(0.454, 0.519)
	α_2	0.3	0.305	0.005	1.843	(0.273, 0.337)
	γ_0	1	0.976	-0.023	0.609	(0.943, 1.009)
	γ_1	1.5	1.533	0.033	1.177	(1.504, 1.563)
	γ_2	-0.2	-0.179	0.020	1.463	(-0.210, -0.148)
	Σ_{22}	16	16.024	0.024	81.311	(15.988, 16.060)
3	α_0	-2	-2.001	0.001	3.859	(-2.013, -1.989)
	α_1	0.5	0.499	-0.004	2.089	(-0.487, 0.511)
	α_2	0.3	0.302	-0.02	1.791	(0.290, 0.314)
	γ_0	1	1.008	0.008	0.572	(0.997, 1.020)
	γ_1	1.5	1.490	-0.009	1.040	(1.479, 1.501)
	γ_2	-0.2	-0.205	-0.005	1.457	(-0.216, -0.194)
	Σ_{11}	4	4.627	0.627	45.457	(3.513, 4.741)
	Σ_{12}	6	6.514	0.514	34.452	(5.324, 6.705)
	Σ_{22}	16	14.293	-1.706	76.523	(13.909, 17.677)

with the increase in the correlation coefficient, we see that the correlated model better fits the simulated data than the other one.

When $\rho = 0$, the separate model showed an increasing bias for the variance components (Σ_{11} , and Σ_{22}) but decreased bias when $\rho = 0.50$; thus, we can conclude that while the results from comparing their MSE and bias for the two models are not significantly different, they do encourage the usage of the correlated spatial CZIDW model, especially when the genuine correlation is significant and performs better than the other. As it is known in Table 2, the MSE values are significantly large, so it is suggested that for future studies of the prior distribution matrix generalized Half-t (MGH-t) (Burger 2020) for the variance-covariance matrix of random effects is used as a replacement for the previous Wishart distribution.

7. ANALYZING THE INFECTION DYNAMICS OF CERCOSPORIOSIS IN OLIVE TREES

To analyze the time to infestation of olive trees affected by *Cercospora*, as detailed in Sect. 2, we will examine the various factors that contribute to the rate of infestation and the overall impact on tree health. This analysis will provide insights into the progression of the disease and its implications for olive cultivation. Suppose A denotes the spatial domain of the orchard, which is divided into $L = 20$ regions A_1, \dots, A_L with clear boundaries. Let $T_{i\ell}$, for $i = 1, \dots, 9, \ell = 1, \dots, L$, be the length of time to infestation for the i th tree

Table 3. Comparison of Models 2 and 3

Model	Parameter	True value	$\rho = 0.75$		$\rho = 0.5$		$\rho = 0.25$		$\rho = 0$	
			Bias	MSE	Bias	MSE	Bias	MSE	Bias	MSE
2	α_0	-2	-0.01	3.868	-0.010	3.903	-0.006	3.91	-0.009	4.127
	α_1	0.5	-0.012	2.136	0.002	2.167	0.002	2.141	0. - 0.016	2.149
	α_2	0.3	0.005	1.843	0.006	1.807	0.007	1.835	-0.010	1.682
	γ_0	1	-0.23	0.606	-0.014	0.600	-0.013	0.618	-0.009	0.575
	γ_1	1.5	0.033	1.177	0.032	1.167	0.026	1.141	0.203	1.264
	γ_2	-0.2	0.020	1.463	-0.004	1.482	0.003	1.485	0.271	1.316
	Criterion		DIC	p_D	DIC	p_D	DIC	p_D	DIC	p_D
3	α_0	-2	-6.327	0.270	-5.732	0.567	-4.629	1.119	-4.3001	1.283
	α_1	0.5	-0.001	3.859	0.008	3.824	0.008	3.812	0.008	3.821
	α_2	0.3	-0.004	2.809	-0.003	2.088	-0.001	2.111	-0.004	2.114
	γ_0	1	0.002	1.791	0.005	1.786	0.003	1.785	0.002	1.769
	γ_1	1.5	0.008	0.572	-0.012	0.567	-0.008	0.568	-0.008	0.570
	γ_2	-0.2	-0.009	1.040	0.006	1.079	0.011	1.506	0.013	1.093
			-0.005	1.457	-0.004	1.460	-0.005	1.459	0.04	1.462
	Criterion		DIC	p_D	DIC	p_D	DIC	p_D	DIC	p_D
			-7.51	0.141	-7.220	0.176	-7.163	0.147	-7.105	0.118

Table 4. Comparison of estimation of variance matrix and covariance of random effects

Model	Parameter	$\rho = 0$			$\rho = 0.25$			$\rho = 0.5$			$\rho = 0.75$		
		True	Bias	MSE	True	Bias	MSE	True	Bias	MSE	True	Bias	MSE
2	Σ_{11}	4	-0.001	53.31	1	-0.005	3.120	4	-0.003	0.269	4	-0.007	48.890
	Σ_{12}	0	-	-	0.5	-	-	4	-	-	6	-	-
	Σ_{22}	16	-0.005	133.32	4	-0.013	7.146	16	0.269	96.368	16	0.024	81.311
3	Σ_{11}	4	0.001	52.95	1	-3.009	78.248	4	0.627	46.781	4	0.331	45.457
	Σ_{12}	0	-0.019	91.412	0.5	-0.493	83.957	4	0.014	47.137	6	-0.251	34.452
	Σ_{22}	16	0.020	133.61	4	-0.022	53.775	16	1.706	92.426	16	0.590	76.523

in the ℓ th square segment which are zero-inflated discrete variables. To fit $T_{i\ell} \mid X, Z \sim CZIDW(\pi_{i\ell}(Z), q_{i\ell}(X), \beta)$ model to the data, first, it is necessary to build design matrixes. According to 1, the $(4 \times n)$ -dimensional design matrix is $X = (1, X_1, \dots, X_3)$, including the covariates “Age” X_1 , “Type” X_2 , and “Height” X_3 . To model count values from a discrete Weibull distribution, the complementary log–log link function

$$\log(-\log(q_{i\ell})) = \alpha_0 + \alpha_1 x_{1i\ell} + \alpha_2 x_{2i\ell} + \alpha_3 x_{3i\ell} + \phi_{2\ell}, \quad (20)$$

is used for the parameter $\mathbf{q} = (q_{i\ell})$. Also, we model the mixing parameter $\boldsymbol{\pi} = (\pi_{i\ell})$, which is the probability of being zero, through the logit link function

$$\text{logit}(\pi_{i\ell}) = \gamma_0 + \gamma_1 z_{1i\ell} + \gamma_2 z_{2i\ell} + \gamma_3 z_{3i\ell} + \phi_{1\ell}, \quad (21)$$

and the design matrix $Z = (1, Z_1, \dots, Z_3)$, with the same covariates with X is considered to model nonzero count values. In the models (20) and (21), $\gamma_{i\ell} = (\gamma_{0i\ell}, \dots, \gamma_{3i\ell})$ and $\alpha_{i\ell} = (\alpha_{0i\ell}, \dots, \alpha_{3i\ell})$ denote the vectors of regression coefficients, $\phi_1 = (\phi_{11}, \dots, \phi_{19})$ and $\phi_2 = (\phi_{21}, \dots, \phi_{29})$ include the spatial correlation of the response variable. We follow the same setup described in Sect. 6, and we assigned weakly informative normal priors to the regression coefficients, a BICAR prior for $\phi_\ell = (\phi_{1\ell}, \phi_{2\ell})^T$, $\ell = 1, \dots, L$, with an inverse Wishart prior for the spatial covariance. We ran the MCMC algorithm for 50,000 iterations with a burn-in of 10,000 and a thin of 5. Model diagnostics, such as the trace plots, suggest that the MCMC chains achieved convergence. The convergence of the chain is determined for each parameter by drawing the effect diagram, i.e., the values of the parameters during the execution time of the chain. Table 5 gives the Bayes estimates and 95% equal-tail credible intervals of the regression coefficients for three models: the CZIDW model, the spatial CZIDW model with independent random effects, and the spatial CZIDW model with spatial correlated random effects.

Table 5 presents a comprehensive analysis of the dynamics of Cercospora disease infestation in olive trees, utilizing Bayesian estimates derived from binomial and count data models.. The Bayes estimates for the binomial component provide intriguing insights into the dynamics of the probability of Cercospora disease infestation within a week. Analyzing the intercepts across various models reveals that all models exhibit positive values, with the CZIDW model estimating approximately 3.019. In contrast, the spatial CZIDW models present slightly elevated estimates, with values of 3.142 for “Independent random effects” and 3.2249 for correlated random effects. This suggests that the spatial models indicate a higher baseline probability of Cercospora disease infestation occurring within a week. Examining the “Age of olive trees” parameter estimates (γ_1) yields mixed results. The CZIDW model indicates a negligible effect of 0.001, while the spatial models, particularly the correlated random effects model, reveal a negative relationship of -0.031 . This finding implies that as trees age, their susceptibility to disease infestation within a week may actually decrease. Regarding the “Type of olive trees” variable (γ_2), the CZIDW model presents a positive estimate of 0.058. However, the spatial models show lower or even negative estimates, especially in the correlated random effects model, which reports -0.028 . This variation suggests that the impact of type on disease susceptibility can differ

Table 5. Bayes estimates and 95% credible intervals for the parameters of three models

Model component	Variable	Parameter	CZIDW	Spatial CZIDW Independent r.e.	Correlated r.e.
Zero component	Intercept	γ_0	3.019 (3.007, 3.030)	3.142 (3.123, 3.160)	3.2249 (3.213, 3.235)
	Age	γ_1	0.001 (-0.010, -0.0127)	-0.006 (-0.025, -0.012,)	-0.031 (-0.043, -0.020)
	Type	γ_2	0.058 (-0.0464, 0.0695)	0.026 (-0.008, 0.045)	-0.028 (-0.038, -0.017)
	Height	γ_3	-0.046 (-0.058, -0.0342)	-0.040 (-0.056, -0.023)	-0.009 (-0.020, -0.002)
	Intercept	α_0	-1.563 (-1.577, -1.553)	-1.558 (-1.575, -1.541)	-1.490 (-1.501, -1.479)
Count component	Age	α_1	-0.002 (-0.013, 0.008)	-0.0167 (-0.033, 0.000)	-0.05 (-0.016, 0.006)
	Type	α_2	0.039 (0.028, 0.050)	-0.007 (-0.025, 0.009)	0.005 (-0.005, 0.015)
	Height	α_3	-0.012 (-0.023, -0.001)	0.020 (0.000, 0.0398)	0.004 (-0.011, 0.013)
Variance components	$Var(\phi_1)$	Σ_{11}	-	0.066 (0.063, 0.069)	0.064 (0.061, 0.066)
	$Cov(\phi_1, \phi_2)$	Σ_{12}	-	-	-0.081 (-0.090, -0.072)
	$Var(\phi_2)$	Σ_{22}	-	1.058 (0.999, 1.117)	2.450 (2.364, 2.535)

significantly depending on the model employed. Lastly, the “Height of olive trees” variable (γ_3) consistently demonstrates negative estimates across all models. This trend indicates that greater height values are associated with a reduced probability of Cercosporiose disease infestation occurring within a week. In the second part, we investigate nonzero count data in conjunction with selected covariates. Our results consistently demonstrate that the intercepts across all models are negative. Notably, our proposed model, the correlated random effects model, exhibits the least negative intercept at -1.490 , which suggests a lower baseline count. Moreover, the “Age of olive trees” parameter (α_1) is negative in all models, indicating a potential decrease in the average time to infestation as trees mature. This finding implies that older trees may experience a reduced duration before becoming infested. The “Type of olive trees” variable presents a more complex narrative; it (α_2) yields a positive estimate in the CZIDW model 0.039 while producing a negative estimate in the independent random effects model -0.007 . Our sample data reveal a multifaceted relationship between tree type and disease susceptibility. Specifically, a change in tree type is positively correlated with the nonzero count model, suggesting that certain tree types may have a longer average time before an olive tree becomes infected. Conversely, the Bernoulli part provides an alternative perspective, indicating a negative coefficient. In our proposed model (spatial correlated random effects model), for each unit change in tree type, the probability of disease occurrence within the first week decreases. This finding underscores the significant influence of tree type on disease probability. The effects of “Height of olive trees” (α_3) also exhibit variability across the models. The independent random effects model indicates a positive estimate 0.020 , while other models report estimates that are either close to zero or negative. For instance, in our proposed model—the spatial model with correlated random effects—the posterior estimate for this parameter is 0.004 , suggesting that a one-unit increase in tree height is associated with an increase of 0.004 in the average time to infection by Cercosporiosis. In summary, our findings indicate that in our proposed model, younger and taller trees tend to exhibit a longer average survival time, with a lifespan of at least one week. This highlights the critical importance of considering both age and height in understanding the dynamics of tree health and disease susceptibility. The variance of the random effect ϕ_1 , denoted as Σ_{11} , is estimated to be approximately 0.066 when considering independent random effects. The corresponding 95% credible interval ranges from 0.063 to 0.069 . In contrast, for the spatially correlated random effects model, the variance is slightly lower at 0.064 , with a credible interval of 0.061 to 0.066 . This suggests that there is some variability in the random effect ϕ_1 , which captures the spatial variation present in the data. For the random effect ϕ_2 , represented as Σ_{22} , the variance is estimated at 1.058 under the independent random effects model, with a credible interval of 0.999 to 1.117 . However, in the spatially correlated random effects model, the variance increases significantly to 2.450 , with a credible interval of 2.364 to 2.535 . This notable difference indicates that there is greater variability in the response variable when accounting for correlated random effects. When comparing the variance components Σ_{11} and Σ_{22} in our proposed model, we gain insights into the variability of the random effects. Specifically, ϕ_2 exhibits greater variability than ϕ_1 , suggesting that there is lower between-segment-group variability in the likelihood of contracting Cercosporiosis disease within a week compared to the average time it takes for an olive tree to become infected. Additionally, our proposed model considers the correlation

Table 6. Comparison of three models

Model	DIC	p_D	WAIC
1	-7.330	0.669	-6.844
2	-7.497	0.752	-6.906
3	-8.395	1.202	-7.038

Table 7. Comparison criteria of spatial models with correlated random effects

Criterion	CZIP	CZINB	CZIDW
DIC	-2.124	-6.351	-8.395
p_D	0.076	0.832	1.202
WAIC	-4.127	-6.359	-7.038

among random effects. Σ_{12} , the covariance between the random effects ϕ_1 and ϕ_2 , is estimated to be -0.081 in the spatially correlated random effects model, with a 95%credible interval ranging from -0.090 to -0.072 . This negative covariance indicates that as one random effect increases, the other tends to decrease, highlighting an inverse relationship between spatial variation of likelihood of contracting Cercosporiosis disease within a week and spatial variation of the average time it takes for an olive tree to become infected. In both spatial cases, the credible intervals of the covariance matrix components were relatively narrow and bounded away from zero, suggesting that the variance components were well identified. Furthermore, the 95% credible interval emphasizes the importance of using the bivariate model and its appropriateness.

The results presented in Table 6 provide a comparative analysis of three models based on two criteria: deviance information criterion (DIC) and widely applicable information criterion (WAIC), along with the effective number of parameters (p_D). There is clear evidence that using a BICAR prior in a spatial model has enhanced the fit of a spatial model when random effects are correlated instead of considering using an intrinsic CAR model. The DIC value for a spatial regression model with spatially correlated random effects -8.395 ($p_D = -1.202$) which is compared to spatial model with independent random effects is reduced. According to DIC and p_D , it also offers better fits than the CZIDW model. The superiority of spatial regression models over classical regression models can be expected due to the spatial correlation of data. The classical model has the worst fit with -7.330 ($p_D = 0.669$) among all considered models. WAIC is another criterion for model comparison, similar to DIC, that accounts for both fit and complexity. Lower WAIC values indicate a better model. In this case, Model 3 again performs the best with a WAIC of -7.038 , followed by Model 2 at -6.906 , and Model 1 with the least favorable WAIC of -6.844 .

Table 7 presents a comparative analysis of the goodness of fit for the spatial model that incorporates correlated random effects associated with the censored zero-inflated discrete Weibull (CZIDW) distribution, in contrast to the censored zero-inflated negative binomial (CZINB) and censored zero-inflated Poisson (CZIP) distributions. The analysis, based on the deviance information criterion (DIC), indicates that the proposed zero-inflated model

utilizing the discrete Weibull distribution significantly outperforms the CZIP and CZINB models. Furthermore, as illustrated in Fig. 3, the presence of over-dispersion in the data suggests that both the CZIDW and the CZINB models, which incorporate correlated spatial random effects, exhibit relatively comparable performance. Both models demonstrate lower DIC values than the CZIP model, indicating a better fit. The effective number of parameters (p_D) serves as an indicator of model complexity, with higher values signifying a more intricate model that may provide a closer fit to the data. The CZIDW model, which exhibits the highest p_D value, suggests a greater utilization of parameters to effectively capture the underlying structure of the data. In contrast, the CZIP model, characterized by the lowest p_D value, indicates a simpler model that may inadequately represent the complexity inherent in the data. In addition to DIC, the widely applicable information criterion (WAIC) is also employed for model evaluation, where lower WAIC values denote superior model performance. Consistent with the DIC findings, the CZIDW model achieves the lowest WAIC value, further reinforcing its superiority in fitting the data compared to the other models evaluated. Overall, the results across all three criteria—DIC, p_D , and WAIC—consistently indicate that the CZIDW model outperforms the CZINB and CZIP models. This finding underscores the enhanced fit of the CZIDW model and its efficacy in addressing the specific characteristics of the analyzed data.

7.1. HAZARD AND ODDS RATIOS OF RISK FACTORS

Table 5 shows the Bayes estimates of regression coefficients in both the count and the zero probability models (Bernoulli models). Their interpretation is the change of the response variable for an increase in one unit in the explanatory variable (risk factor). Since covariates in the count model are related to q through the complementary log–log link function and in the Bernoulli model are associated with π through the logit link function, interpreting these coefficients is difficult in such a way that the role of each subgroup of covariates could be clearly defined in the two components of the model, the probability of early Cercorpiose disease infestation, the average “time-to infestation” provided that the tree has survived the disease for at least a week. Therefore, it is necessary to calculate the criteria that can compare the risk of exposure to a variable to not being exposed to it.

The value of odds ratio $OR = p/(1 - p)$ can be 0 to infinity. A value of 1 indicates no association between the risk factor and the outcome. A value greater than 1 indicates a positive association, i.e., the risk factor increases the chance of the event. In contrast, a value less than 1 shows a negative association, i.e., the risk factor decreases the possibility of the event. In logistic regression, the odds ratio through the exponential function of regression coefficients $\frac{\exp(\hat{\gamma})}{1+\exp(\hat{\gamma})}$ is calculated. This criterion allows one to compare the role of different covariates (risk factors) in the probability of zero survival.

The hazard ratio $HR = (\text{Hazard rate in group2}) / (\text{Hazard rate in group 1})$ allows us to compare the two groups' final event occurrence rates. It can assign a specific number to the additional risk an individual in one group is exposed to compared to an individual in another group. With HR, a result of 1 implies that both groups have the same sum of chance, whereas it is not equal to 1, showing that one group bore more opportunity than another. A value less than 1 indicates that the event's occurrence rate is lower in an exposed group than in a

Table 8. Bayes estimates and 95% credible intervals for OR and HR in a spatial CZIDW model with correlated random effects

Risk factor	OR		HR	
	Est.	Credible interval	Est.	Credible interval
Age	0.033	(0.001 , 0.142)	1.643	(1.256 , 1.858)
Type	0.975	(0.853 , 1.237)	2.906	(2.021 , 3.327)
Height	0.911	(0.762 , 1.143)	3.301	(2.431 , 3.782)

non-exposed group. The HR is a helpful measure because it is straightforward to interpret and lets us know immediately whether exposure to a risk factor increases or decreases the hazard rate. Since the link function used in the counting model is the complementary log–log function, the hazard rate ratio equals $\exp(\hat{\alpha})$.

Table 8 shows the Bayes estimates and 95% credible intervals of OR and HR for significant risk factors in the two-part models: Bernoulli model for zero component via logit link function and discrete Weibull model for count component via complementary log–log link function. The chance that the time to infection with *Cercosporiosis* fungus is more than a week in the group of younger trees is 1.64 times more than in the older one. Also, the risk ratio for the shorter trees that have resisted this disease for at least one week is 3.3 times higher than for taller trees with the same characteristics. However, the odds ratio for disease incidence in the first week is 0.911 for taller trees compared to shorter trees, and for older trees, it is 0.033 compared to younger trees.

7.2. ANALYSIS OF TREE'S AGE AT DISEASE TIME

The age of an olive tree significantly impacts its susceptibility to diseases such as *Cercosporiosis* leaf spot, caused by the fungus *Cercospora*. With their developing immune systems and less established root systems, younger trees may initially appear more vulnerable to infections. However, they often exhibit vigorous growth and a stronger immune response, which can enhance their resilience against certain diseases.

In contrast, older, well-established olive trees typically possess more robust defenses, making them better equipped to withstand diseases. Their adaptability is a source of reassurance, as they can become susceptible to infections if they experience stress from environmental factors such as drought, poor soil conditions, or inadequate care. Such stressors can weaken their overall health and increase their vulnerability to pathogens, but their robust defenses provide a buffer against these threats.

Additionally, the age of the tree influences the manifestation of disease symptoms. Older trees may show signs of stress or infection more slowly than younger ones, complicating early detection and management. Furthermore, as trees age, they may accumulate pathogens and experience physiological changes, such as reduced growth rates and alterations in leaf structure, which can further affect their interaction with diseases like *Cercospora*.

Given that the age of each olive tree is a time-dependent random variable and an essential risk factor in susceptibility to diseases, we consider age as a smooth function and fit it by a cubic B-spline. With internal nodes in the first, second, and third quartiles, we estimated

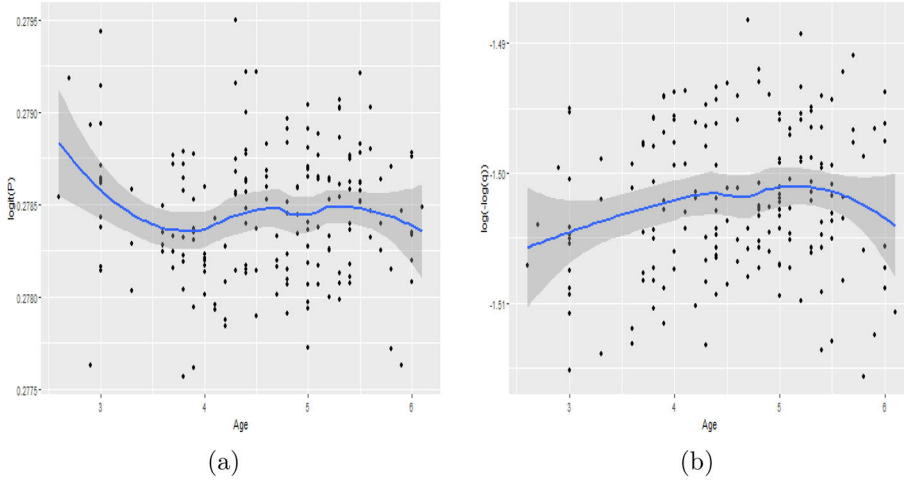


Figure 9. **a** The effect of changes in an olive tree's age overtime on the probability of survival time becoming zero; **b** The impact of changes in the olive tree's age overtime on the duration of time to disease.

the age distribution for women, respectively (4, 4.7, 5.3) year. This reiterates the gravity of our research in understanding the age effect as a risk factor in disease susceptibility. In this section, we intend to consider the age effect as a time-varying covariate on the linear predictor for the Bernoulli component and the discrete Weibull component of the two-part spatial CZIDW model with correlated random effects. As the tree age increases to before four years, the chance of disease in the first week decreases drastically, as shown in Fig. 9a. From the age of 4 onward, this chance changes upward until the age of 5. However, around five years old, the tree stabilizes after a slight decline.

In Fig. 9b, we illustrate the potential impact of changing the age of trees resistant to this disease for at least one week on average until the disease occurs. The figure clearly shows that increasing the age of the trees up to four years increases the resistance time of the tree against this fungus. However, at 4.5 years old, we observe a slight fluctuation (downward–upward) in disease resistance. Finally, as the tree's age increases from 5 years onward, its survival time decreases, underscoring the urgent need to understand and address the potential impact of tree age on disease resistance.

7.3. PREDICTING TIME TO DISEASE OF CERCOSPORIOSIS

In this subsection, we have predicted disease time for high-risk groups for the probability of fungal disease of olive trees in the first week through the CZIDW spatial survival model with correlated random effects. This group includes young trees, type one and shorter, which is shown in Fig. 10a. The bolder parts are the areas where the risk of this disease is predicted to be higher for the high-risk group in the first week than other parts. Also, Fig. 10b shows the estimated number of weeks until disease occurs for a group of “high-risk” trees, provided they are disease-free for at least the first week. The darker parts of the garden are where trees that have resisted this fungus for at least a week have a shorter lifespan.



Figure 10. **a** The probability of contracting the disease in less than a week, **b** The average time until the disease is infected, provided that the tree has resisted the infection for at least one week .

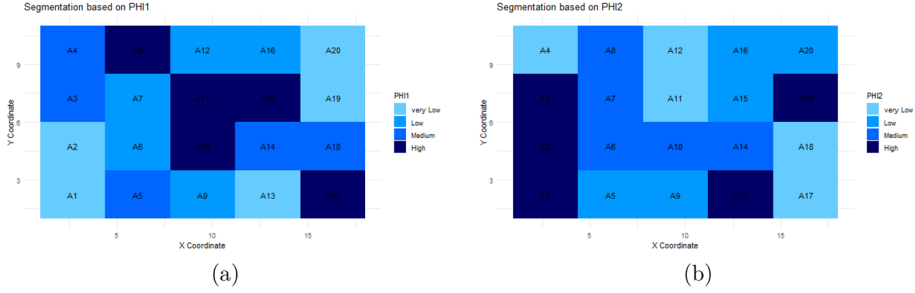


Figure 11. Geographical distribution of model-based predictions of spatial random effects **a** ϕ_1 , **b** ϕ_2 .

Figure 8 shows the geographic distribution maps of model-based predictions of spatial random effects ϕ_1 and ϕ_2 and emphasizes the significance of regional differences of random effects. Dark blue segments generally have higher predicted values ($\phi_r = 0, r = 1, 2$) than others. At the same time, light blue areas have lower values in random effects (i.e., $\phi_r < 0, r = 1, 2$). Because the random effect ϕ_1 is used to enter the spatial correlation between area via Equation (21), according to Fig. 8a, it can be concluded that spatial factors after controlling for other factors in some segments have the most significant spatial effect on the probability of disease in less than a week in this area. Also, in Fig. 8b, the geographical distribution of random effect ϕ_2 is depicted. This random effect enters the spatial correlation between the areas via Equation (20).

8. DISCUSSION AND RESULTS

The choice of an appropriate distribution is paramount in the parametric modeling of data with an excessive number of zeros. The ZIDW distribution, capable of reflecting any dispersion, over-dispersion, equal distribution, or under-dispersion, is a crucial selection. This distribution not only aids in accurate parameter estimation but also enhances data interpretability in survival studies, thereby empowering statisticians and researchers. Therefore, the ZIDW distribution, a generalization of the widespread continuous Weibull distribution in survival studies, was chosen in this study. Our approach was thorough, considering the two-component nature of zero-inflated models and the corresponding two-part spatial model.

We took two distinct approaches to include random spatial effects in the model, ensuring a comprehensive exploration of the topic. By conducting simulation studies, we found higher parameter estimation accuracy in the spatial model with correlated random effects. After examining the evaluation criteria of the existing models on the data of the time to get Cercosporiosis disease caused by *Cercospora* fungus event with inflation in the 0 to 7 days, we found that this model fits better than other censored zero-inflated discrete models like Poisson and negative binomial. Nevertheless, in all these models, spatial variety is considered in the spatial random effects, and the spatial structure of covariates is overlooked. So, the estimation of parameters needs to be more precise. Even if the covariates directly affect the response variable without any observable spatial variability, they can still be impacted by confounding variables that do show spatial variation. This means other factors varying across space might influence the relationship between covariates and the response variable. Accordingly, it would be possible that an unmeasured spatially varying confounder caused the spatial correlation in the random effects. In addition, the random effect approach does not include the geographical risk in the definition of the survival measures and, therefore, has limited interpretability. Thus, methodologies that avoid using spatial random effects would be better in the zero-inflated spatial survival model, especially in the correlation between the covariates and spatial structure.

ACKNOWLEDGEMENTS

The authors thank the editorial team and the anonymous reviewers for their comments, helpful suggestions, and encouragement, which helped improve the final version of this paper. Receiving support from the Center of Excellence in Analysis of Spatio-Temporal Correlated Data at Tarbiat Modares University is acknowledged.

Declaration

Conflict of interest The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

[Received January 2024. Revised February 2025. Accepted February 2025.]

REFERENCES

- Allison PD (1982) Discrete-time methods for the analysis of event histories. *Sociol Methodol* 13:61–98
- Banerjee S, Carlin BP, Gelfand AE (2004) Hierarchical modeling and analysis for spatial data. Chapman and Hall/CRC, Boca Raton, pp 1–95
- Banerjee S, Dey DK (2005) Semiparametric proportional odds models for spatially correlated survival data. *Life-time Data Anal* 11:175–191
- Besag J (1974) Spatial interaction and the statistical analysis of lattice systems (with discussion). *J R Stat Soc B* 36:192–236
- Besag J, York J, Mollié A (1991) Bayesian image restoration with two applications in spatial statistics. *Ann Inst Stat Math* 43:1–59
- Braekers R, Grouwels Y (2016) A semi-parametric Cox's regression model for zero-inflated left-censored time to event data. *Commun Stat Theory Methods* 45:1969–1988

- Burger DA, Schall R, Ferreira JT, Chen DG (2020) Arobust Bayesian mixed effects approach for zero inflated and highly skewed longitudinal count data emanating from the zero inflated discrete weibull distribution. *Stat Med* 39(9):2
- Calsavara VF, Rodrigues AS, Rocha R, Louzada F, Tomazella V, Souza AC, Costa RA, Francisco RP (2019) Zero-adjusted defective regression models for modeling lifetime data. *J Appl Stat* 46:2434–2459
- Carlin BP, Banerjee S (2003) Hierarchical multivariate CAR models for spatio-temporally correlated survival data (with discussion). In: Bernardo JM, Bayarri MJ, Berger JO, Dawid AP, Heckerman D, Smith AFM, West M (eds) *Bayesian statistics*, vol 7. Oxford University Press, Oxford, pp 44–63
- Chanialidis C, Evers L, Neocleous T, Nobile A (2017) Efficient Bayesian inference for COM-Poisson regression models. *Stat Comput*. <https://doi.org/10.1007/s11222-017-9750-x>
- Chen J, Chen Z (2008) Extended Bayesian information criteria for model selection with large model spaces. *Biometrika* 95(3):759–771. <https://doi.org/10.1093/biomet/asn034>
- Cox DR (1972) Regression models and life-tables (with discussion). *J R Stat Soc B* 34:187
- Da Silva MF, Ferrari SL, Cribari-Neto F (2008) Improved likelihood inference for the shape parameter in Weibull regression. *J Stat Comput Simul* 78(9):789–811
- Darmofal D (2006) Spatial econometrics and political science. <https://api.semanticscholar.org/CorpusID:34295207>
- Eilers PHC, Marx B (1996) Flexible smoothing with B-splines and penalties. *Stat Sci* 11(2):89–121
- Gamerman D, Lopes HF (2006) Markov chain Monte Carlo stochastic simulation for Bayesian inference. Taylor and Francis, Boca Raton
- Gelfand AE, Vounatsou P (2003) Proper multivariate conditional autoregressive models for spatial data analysis. *Biostatistics* 4:11–25
- Gilks WR, Richardson S, Spiegelhalter D (1995) Markov chain Monte Carlo in practice. Taylor & Francis, Milton Park
- Haselimaashhadi H, Vinciotti V, Yu K (2016) A new Bayesian regression model for counts in medicine. *arXiv: Methodology*
- Hastings WK (1970) Monte Carlo sampling methods using Markov chains and their applications. *Biometrika* 57:97–109
- Hennerfeind A, Brezger A, Fahrmeir L (2006) Geoadditive survival models. *J Am Stat Assoc* 101:1065–1075
- Hoff PD (2009) A first course in Bayesian statistical methods. Springer, New York
- Jenkins SP (2005) Survival analysis. Unpublished Manuscript, Institute for Social and Economic Research, University of Essex, Colchester
- Kalktawi HS (2017) Discrete Weibull regression model for count data. Ph.D. thesis; Brunel University London
- Kaplan EL, Meier P (1958) Nonparametric estimation from incomplete observations'. *J Am Stat Assoc* 53(282):457–481
- Khan MSA, Khalique A, Abouammoth AM (1989) On estimating parameters in a discrete Weibull distribution. *IEEE Trans Reliab* 38(3):348–350
- Klakattawi HS, Vinciotti V, Yu K (2018) A simple and adaptive dispersion regression model for count data. *Entropy* 2:142. <https://doi.org/10.3390/e20020142>
- Lambert D (1992) Zero-inflated Poisson regression, with an application to defects in manufacturing. *Technometrics* 34(1):1–14
- Lawless JF (2011) Statistical models and methods for lifetime data. Wiley, Hoboken, p 362
- Lawson AB, Banerjee S, Haining RP, Ugarte MD (2016) Handbook of spatial epidemiology. CRC Press, Boca Raton
- Lee Y, Alam MM, Noh M, Rönnegråd L, Skarin A (2016) Spatial modeling of data with excessive zeros applied to reindeer Pellet-group counts. *J Ecol* 6(19):7047–7056
- Loquiha O, Hens N, Chavane L, Temmerman M, Osman N, Faes C, Aerts M (2018) Mapping maternal mortality rate via spatial zero-inflated models for count data: a case study of facility-based maternal deaths from Mozambique. *PLoS ONE* 13(11):e0202186

- Louzada F, Moreira FF, de Oliveira MR (2018) A zero-inflated non default rate regression model for credit scoring data. *Common Commun Stat Theory Methods* 47:3002–3021
- Mardia KV (1988) Multi-dimensional multivariate Gaussian Markov random fields with application to image processing. *J Multivar Anal* 24(2):265–284
- Metropolis NR (1953) Equation of state calculations by fast computing machines. *J Chem Phys* 21(6):1087–1092
- Moran PAP (1950) Notes on continuous stochastic phenomena. *Biometrika* 37(1/2):17–23. <https://doi.org/10.2307/2332142>
- Moriña D, Puig P, Navarro A (2022) Analysis of zero-inflated dichotomous variables from a Bayesian perspective: application to occupational health. *BMC Med Res Methodol* 22(1):210
- Motarjem K, Mohammadzadeh M, Abyar A (2020) Geostatistical survival model with Gaussian random effect. *Stat Pap* 61(1):85–107
- Nakagawa T, Osaki S (1975) The discrete Weibull distribution. *IEEE Transactions on Reliability* 24(5):300–301
- Neelon BH, Ghosh P, Loebs PF (2015) A spatial Poisson hurdle model for exploring geographic variation in emergency department visits. *J R Stat Soc Ser A* 176:389–413
- Neelon B, Zhu L, Neelon SE (2015) Bayesian two-part spatial models for semicontinuous data with application to emergency department expenditures. *Biostatistics* 16(3):465–79
- Nyandwi E, Osei FB, Amer S, Veldkamp A (2020) Modeling schistosomiasis spatial risk dynamics over time in Rwanda using zero-inflated Poisson regression. *Sci Rep* 10:1–9
- Oliveira MR, Moreira F, Louzada F (2017) The zero-inflated promotion cure rate model applied to financial data on time-to-default. *Cogent Econ Finance* 5:1395950
- Rideout M, Hinde J, Demetrio CGB (2001) A score test for testing a zero-inflated Poisson regression model against zero-inflated negative binomial alternatives. *Biometrics* 57:219–223
- Schnell P, Bandyopadhyay D, Reich BJ, Nunn M (2015) A marginal cure rate proportional hazards model for spatial survival data. *J R Stat Soc: Ser C: Appl Stat* 64(4):673–691
- Souza H, Louzada F, Ramos LM, Oliveira Júnior MR, Perdoná GS (2022) A Bayesian approach for the zero-inflated cure model: an application in a Brazilian invasive cervical cancer database. *J Appl Stat* 49(12):3178–3194
- Souza H, Louzada F, Perdoná GS (2022) The log-normal zero-inflated cure regression model for labor time in an African obstetric population. *J Appl Stat* 49(9):2416–2429
- Spiegelhalter DJ, Best NG, Van der Linde A (2002) Bayesian measures of model complexity and fit (with discussion). *J Roy Stat Soc B* 64:583–639
- Sturman MC (1999) Multiple approaches to analyzing count data in studies of individual differences: the propensity for Type I errors, illustrated with the case of absenteeism prediction. *Educ Psychol Meas* 59(3):414–30
- Tibshirani R (1996) Regression shrinkage and selection via the Lasso. *J R Stat Soc B* 58:267–288
- Tin A (2008) Modeling zero-inflated count data with underdispersion and overdispersion, research foundation for mental hygiene. *SAS Global Forum, Statistics and Data Analysis*, 372
- Ver Hoef JM, Jansen JK (2007) Spacetime zero-inflated count models of harbor seals. *Environmetrics* 18:697–712
- Wang Z, Shuangge Ma, Wang CY (2015) Variable selection for zero-inflated and overdispersed data with application to health care demand in Germany. *Biom J* 57(5):867–884
- Watanabe S (2010) Asymptotic equivalence of Bayes cross validation and widely applicable information criterion in singular learning theory. *J Mach Learn Res* 11:3571–3594
- Zhu H, DeSantis SM, Luo S (2018) Joint modeling of longitudinal zero-inflated count and time-to-event data: a Bayesian perspective. *Stat Methods Med Res* 27:1258–1270

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.