

Modelling spatial point patterns of Asian Giant Hornets' occurrences

Baihui Che

Nanyang Technological University, 50 Nanyang Ave, Singapore 639798

bche002@e.ntu.edu.sg

Abstract. This study explores statistical methods for modelling the spatial patterns of Asian giant hornets using sighting records in the Republic of Korea from 2008 to 2013. The focus is on simulating the spatial distribution of hornets, with a key aspect being the statistical inference of their intensity. Gaussian kernel-smoothing estimation is utilized to model the hornets' intensity based on sighting records, which also transforms the occurrences of honeybees, a primary prey of the hornets, into a covariate. Results show a significant dependency of hornets' intensity on honeybees' intensity, with a calculated probability that a hornet sighting location has higher honeybee intensity than a random location. By this finding, the parametric modelling of the hornets' spatial intensity is applied with the covariate of honeybees in each year, with the basic inhomogeneous Poisson point process along with the log-linear model. The model is refined for each year using backward stepwise selection based on the Akaike Information Criteria. Model validation confirms the Poisson process assumption and shows promising results with raw residuals against honeybee intensity. The analysis demonstrates that the spatial pattern modelling method employed is both sensible and valid.

Keywords: Poisson point process, Spatial pattern analysis, Ecological modelling

1. Introduction

1.1. Background

The Asian giant hornet (AGH), which is scientifically called *Vespa mandarinia* Smith, 1852, is native to Asia and seen as an invasive species on American continents in recent years. In September 2019, a colony of AGH was discovered in Vancouver, Canada. Although the nest was quickly eradicated, several confirmed sightings of the pest have occurred in neighbouring Washington State since that moment [1]. A reason for people to be intensively repugnant to the giant hornets is that they are well known for preying on honeybees as their food and a source of protein for their larvae, and destructing colonies or hives of the prey, leading to the financial loss of beekeepers [2]. The other reason for aversion comes from the fatal attack of AGH to human beings. In 2020, a report in Japan claims that the giant hornet can spray venom into people's eyes and cause nearly permanent injury [3]. In a more severe incident in 2013, attacks by AGH in Shaanxi, China, resulted in 42 fatalities and over 1,600 injuries [4].

Vespa mandarinia, unlike typical bees and hornets that construct hives on trees or walls, tends to establish its nests near tree roots or in existing tunnels created by other animals [5], effectively concealing their presence. Moreover, it seems that people may easily misclassify them as other wasps

or hornets like *V. velutinia* [6]. Therefore, the sighting record is valuable to some extent, which enables the behaviour of AGH such as spread pattern and static distribution from the geographic perspective to be studied, for the goal of informing people more about this dangerous species and contributing to the preservation against it. In the paper, the spatial point pattern analysis in 2 dimensions is utilized to investigate the geographic occurrence record of this ‘murder hornet’. The objective is to statistically infer the properties of the spatial distribution of sighting locations and to implement point process models to decipher the behaviour of the insects. By simulating occurrences over a certain period and examining the evolution of point patterns, the study aims to derive significant insights into the dispersal behaviour of AGH. These findings are anticipated to contribute to the broader research on potentially invasive insect species in various continents [7].

1.2. Literature Review

The study intersects spatial statistics, point processes, entomology, and ecology. It explores spatial point patterns—locations in 2 or 3 dimensions within a region, generated by stochastic mechanisms [8]. Its applications span various fields like astronomical structures [9], ecological patterns in plant communities [10], particle distribution [11], and epidemiology [12]. In entomology, spatial point pattern analysis, like the study of water stride larvae distribution [13], employs models like the uniform Poisson process and Strauss point process for species data analysis. These tools also apply to human behaviour studies, such as crime map analysis [14], focusing on crime dependency and dispersion patterns.

In the field of entomology, particularly concerning the dispersal patterns of species, there are limited studies on *Vespa mandarinia* in the context of the Republic of Korea, with more focus on species like *V. velutina* [15] and *V. simillima* [16]. Studies often cover ecological aspects like taxonomy [15], morphology [17], and urbanization impacts [15]. The Asian giant hornet’s global spread is increasingly studied using ecological niche modelling (ENM) for distribution prediction with environmental and human factors [2], climate change [18], and honey production and species diversity [19]. ENM, or species niche modelling [20], uses algorithms with climatic and environmental data [21] employing techniques like regression and machine learning (SVM, ANN). It requires extensive biological data and may be limited by data availability. In contrast, spatial point patterns analysis relies less on external data, focusing on the statistical properties of objects, emphasizing pattern analysis over prediction.

In the paper, spatial statistics is used to try to analyze the 2-dimension location of AGH, and this is de facto not a typical case for distribution investigation such as the gorilla nests data by the Wildlife Conservation Society Takamanda-Mone Landscape Project [22] because the observations are Asian giant hornets, the insects which may appear in some places but are unable to be seen by people for some non-technical reasons. Unlike hives, colonies, and nests that are typically stationary and easier for observers to locate, the distribution of this species is not as straightforward, given its propensity to move to different locations, making detection more challenging. In other words, it is likely that some activities cannot be found, which may influence the statistical inference results. Some research to predict the potential spread of AGH ignores this problem [2]. However, though this is unavoidable, to some degree it is acceptable since statistics can take the advantage of itself to draw a meaningful conclusion from imperfect samples. Therefore, it is reasonable to use spatial point pattern modelling to pursue the study of this topic.

1.3. Materials and Methods

The data studied is obtained from GBIF [23], it is stored in the National Institute of Ecology in the Republic of Korea, which recorded the occurrences by human observation with latitudes, longitudes, and exact dates of sightings from 2008 to 2013 about *Vespa mandarinia*, the hornet, and *Apis mellifera* Linnaeus, the honeybee. The first reason to choose the data in South Korea is that AGH is a native hornet there [24], so the history of living and spread of it shall be relatively longer than most other regions where it is not local. The second reason is that the country contains more records and different locations than any other countries, which is better for the operation of spatial statistics study. The third factor is

the six-decimal precision of coordinates, essential for the uniqueness of each sighting and preventing data duplication that lesser decimal places might cause.

The duplicated points are useless in spatial pattern modelling because they may impede some techniques of the methodology [25]. Thereby, when using the point pattern process in the study, duplicated points would be deleted. Unlike the data of South Korea, the one of Japan [26] contain unclear coordinates which leads to many coincident points and then may make the modelling procedure invalid, though the records are much more abundant.

The research utilizes statistical tools like the intensity function, pair correlation function, and K-function for spatial pattern analysis, particularly in constructing Poisson point processes. It examines pattern distributions across the continent over time and in smaller areas. The distribution in the continent during years and patterns within a smaller observation window are explored. In the study, the numerical covariate is included, and spatial analysis is conducted using R Statistical Software with the spatstat toolbox [27].

2. Complete Spatial Pattern

2.1. Choice of Observation Window

Here the study first considers the full spatial pattern formed by the data from 2008 to 2013, as it aims to do the exploratory data analysis from the whole distribution primarily. However, the cumulative patterns through the years should still be displayed to facilitate the subsequent contemplation in the study. The series of plots as followed in Figure 1 and 2 shows the enlarging range of human sighting patterns of AGH in the continent of the Republic of Korea. From the plots, it can be found that if the image in 2008 is assumed to be the start of the distribution pattern, then visually the spread of sightings began from the top, middle, and bottom and continued to cover the whole area in the largest continent in the country. In this section, the last pattern is studied.

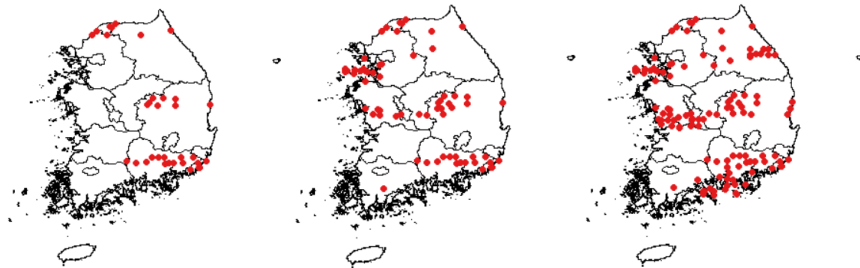


Figure 1. Cumulative occurrences maps from 2008 to 2010 (Left to Right).



Figure 2. Cumulative occurrences maps from 2011 to 2013 (Left to Right).

In the interest of investigating the point pattern in the map, it is necessary to establish an observation window. An observation window constrains the scope of research, which statistical properties and fitted models are concluded based on. In general, the instruction that better statistical results are yielded with larger window holds [28]. As is shown in Figure 3, it is found in the southern part of the continent, there are too many complicated and irregular shapes on the border. If this part is included in the analysis, even

though the result and the predictor or simulation are generated, the conclusion may not be reasonable since the area of sea between the land where sighting cannot arise, is involved in the research without control. This is the area that should not be examined in the study. Hence, the points located in these regions are not counted here, and what is of interest are the points inland. A smaller window may contain misleading information and make researchers put forward conclusions with significant error and bias. Therefore, an observation window as broad as possible is the goal, but for simplicity, the polygonal or irregular shape is not considered, since this is computation-consuming and inefficient, so the standard rectangular shape is chosen as the observation window W which should guarantee the area within the window is inland.

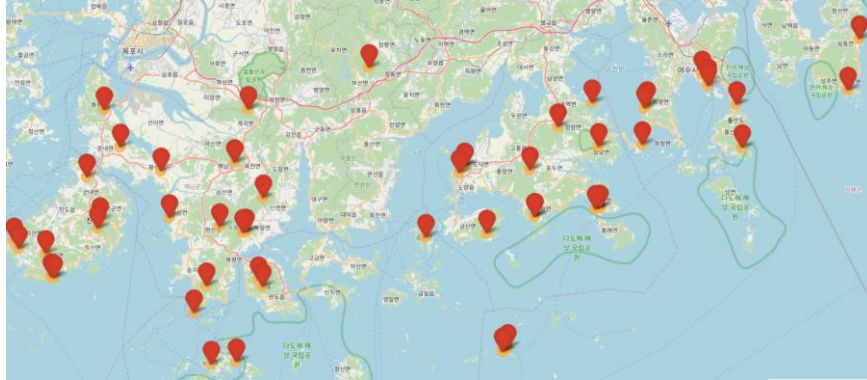


Figure 3. Map of continental boundary against sea with occurrences marks.

2.2. Intensity Estimation

Intensity in spatial point patterns measures the expected number of points in a unit area u . Suppose there is a point process Q , then the expected number of points of Q falling in a region S is given by

$$\mathbb{E}[n(Q \cap S)] = \int_S \lambda(u) du,$$

where $\lambda(u)$ is the intensity function and u indicates a unit area, and the expectation of number is equal to $\lambda|S|$ if the intensity is constant, where $|S|$ is the area of region S .

To explore the intensity function for modelling, we should find whether the pattern is inhomogeneous. Technically, one simple statistical test called the quadrat counting test can be used under the assumption of homogeneity of intensity [29]. Here the hypotheses are set as

H_0 : The intensity is homogeneous,

H_1 : The intensity is inhomogeneous in some unspecified mode.

Then pieces of quadrat are separated as multiple identical m rectangles in the observation window here (See Figure 4). The number of points in a quadrat is reflected by colours. Set each piece contains number of points n_i , then the expected count in each square is equal to $e_i = \frac{\sum n_i}{m}$. The test statistic would be

$$X^2 = \sum_i \frac{(n_i - e_i)^2}{e_i}.$$

Under the null hypothesis, the distribution of X^2 is approximately a χ^2 distribution with $m-1$ degrees of freedom. In the window W , the p-value of the test is 2.55×10^{-7} , which implies that the null hypothesis is rejected under 99% confidence.

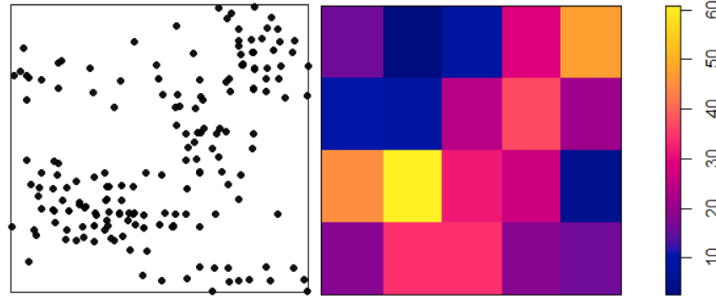


Figure 4. The pattern within observation window W (left) and quadrat counting 4×5 plot (right).

Thereby, the spatial pattern is inhomogeneous, so it seems that the points (Asian giant hornets) have preference for spatial locations. When the homogeneity is assumed, the spatial pattern satisfies complete spatial randomness (CSR) which is simulated by the homogenous Poisson point process. For the research of the spatial pattern, the inhomogeneous Poisson process can be simply assumed.

A modelling method called kernel estimation can estimate the intensity. Though the observation is discrete, the method estimates in a continuous way. Before introducing this approach, an important problem, the edge effect [30], usually occurs when computing empirical characteristics of point patterns. Imagine there is a Poisson point process displayed in Figure 5, and an average number of neighbours within a distance r originated in a typical random point in an observation window is of interest. Points recorded only within the observation window; situations outside are often missed, and not addressing this leads to “uncorrected” empirical results.

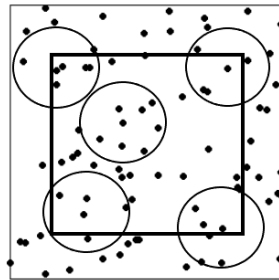


Figure 5. Homogenous Poisson Point Process with $\lambda = 100$ in a unit area.

Basic definitions for intensity estimation in point patterns are introduced, considering summary characteristics. Despite the seeming contradiction of stationarity in this context [28], it is assumed for applying summary characteristic formulas. This case is distinct: observations primarily encompass inland South Korea, rendering the point pattern representative. Unlike usual cases where non-stationary patterns risk misinterpretation due to unrepresentative samples, this study’s cautious use of summary characteristics mitigates such risks. Ripley’s K-function [31] is defined as the expected number of r -neighbours of a typical point of a point process X divided by the intensity. Its empirical form is given by \hat{K} , which is given by

$$\hat{K}(r) = \frac{1}{\hat{\lambda}(x)^2 |W|} \sum_i \sum_{j \neq i} e_K(x_i, x_j) I\{\|x_i - x_j\| \leq r\},$$

where x is the point pattern, $e_K(x_i, x_j)$ is the edge correction weight, and $\hat{\lambda}(x)^2 = \frac{n(x)(n(x)-1)}{|W|^2}$ is an unbiased estimator of the square of intensity in CSR.

Following $\hat{K}(r)$, the pair correlation function [28], which measures the number of pairs of objects by observations in the pattern that is distant from r units apart, divided by the expected number that is seen if the points satisfy CSR, is defined as

$$g(r) = \frac{l}{2\pi r} \frac{dK(r)}{dr}.$$

One method of solving the edge effect called the isotropic correction method [32] is applied in 3.3.4. Suppose there are two points x_1 and x_2 with distance $r = \|x_1 - x_2\|$ from a point process X within a window W , so x_2 lies in the circle $b(x_1, r)$ centred at x_1 with radius r . Then the probability that x_2 lies inside W is demonstrated as the fraction of length of the circle lying within W

$$\Pr(x_1, r) = \frac{l(W \cap \partial b(x_1, r))}{2\pi r},$$

where $\partial b(x_1, r) = \{u: \|x_1 - u\| = r\}$. The K-function adjusted called Horvitz-Thompson estimator by isotropic correction method is then given by

$$\hat{K}_{iso}(r) = \frac{l}{\lambda n} \sum_i \sum_{j \neq i} \frac{l}{\Pr(x_i, \|x_i - x_j\|)} I\{\|x_i - x_j\| \leq r\}.$$

For the uniformly corrected estimators (Bithell, 1990), the true intensity being homogeneous implies the unbiasedness of the estimator, and it gives the study a reason to choose not to adopt it because it has been shown above that the true intensity is inhomogeneous. Hence another method called Diggle's correction [34] is chosen here. This method is expected to have a smaller mean square error with the integral equal to the number of points detected over the window. Set v, x_i as some locations, n as the number of points, and the kernel estimator of the intensity function inside a window W by Diggle's correction is given by

$$\tilde{\lambda}^{(D)}(u) = \sum_i^n \frac{l}{e(x_i)} \kappa(u - x_i),$$

where $e(u) = \int_W \kappa(u - v) dv$ is a correction of bias resulting from edge effects, and $\kappa(\cdot)$ is isotropic Gaussian probability density.

Kernel bandwidth σ , crucial for smoothing the intensity function, is selected through algorithms aimed at minimizing error measures. One notable method is the mean square error (i.e., $E[\{\lambda(S) - \Lambda(S)\}^2]$ where $\Lambda(S) = E[n(X \cap S)]$) cross-validation [35], based on the assumption of a Cox process pattern model. By this procedure, the best σ is found when

$$M(\sigma) = \frac{MSE(\sigma)}{\lambda^2} - g(0)$$

achieves the minimum where $MSE(\sigma)$ is the mean square error at bandwidth σ , λ is the mean intensity, and $g(\cdot)$ is the pair correlation function.

The other one is called the likelihood cross-validation method [36]. The method instead assumes the pattern to be an inhomogeneous Poisson process, which is exactly what the analysis of CSR indicates in the case. By this procedure, the best σ is found when the likelihood cross-validation criterion

$$CV(\sigma) = \sum_i^n \log(\lambda_{-i}(x_i)) - \int_W \lambda(u) du$$

achieves the maximum where $\lambda_{-i}(x_i) = \sum_{j \neq i} \frac{1}{e(x_i)} \kappa(x_j - x_i)$ is the leave-one-out kernel-smoothing estimate of the intensity at x_i with σ , $\lambda(u)$ is the kernel-smoothing estimate of the intensity at location u with σ . The different assumptions of the two methods may cause considerable difference between the optimized σ respectively, but it is not the case here (see Figure 6), and the reason is probably that the actual process is much alike both two models. Note that the left plot is an image of the mean square error cross-validation method, and the right one is an image of the likelihood cross-validation method. Optimized σ is 0.1400 and 0.1508 respectively.

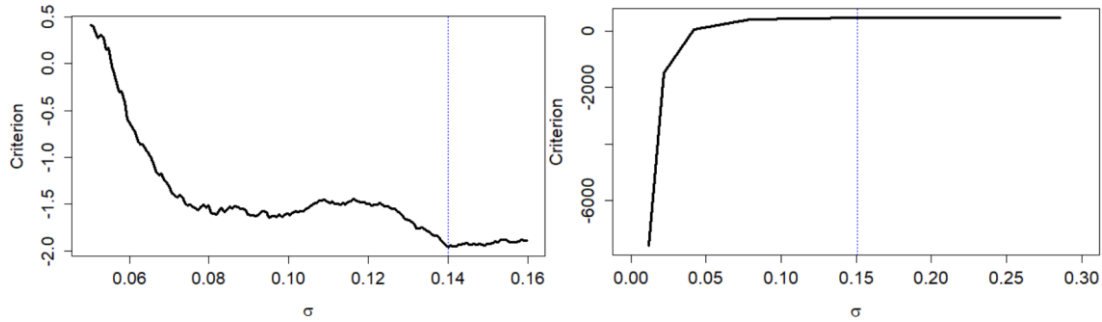


Figure 6. Kernel bandwidth σ the intensity function

After obtaining the σ , the standard error of the intensity can be derived, but only with the assumption that the process is Poisson point process (Daley and Vere-Jones, 1988), as it is presumed at the beginning, so only the full result of likelihood cross-validation is illustrated here. Since the isotropic Gaussian kernel is chosen and this is a 2-dimension case, $\kappa(x, y|\sigma) = \frac{1}{2\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}}$. Given the unbiased estimator of $var(\lambda(u))$ is $\sum_i \frac{1}{e(x_i)^2} \kappa(u - x_i)^2$ [37], the variance takes the form of intensity estimation. Therefore, it can be derived by smoothing the data by the kernel function with σ_2 and times a constant, because

$$\kappa(x, y|\sigma)^2 = \frac{1}{4\pi^2\sigma^2} e^{-\frac{x^2+y^2}{\sigma^2}} = \frac{1}{4\pi\sigma^2} \frac{1}{2\pi(\frac{\sigma}{\sqrt{2}})^2} e^{-\frac{x^2+y^2}{2(\frac{\sigma}{\sqrt{2}})^2}},$$

so $\sigma_2 = \frac{\sigma}{\sqrt{2}}$ and the constant is $\frac{1}{4\pi\sigma^2}$. The result is demonstrated with contour plots, heatmap and a 3D map as followed (See Figure 7). The top left contour plot displays intensity with labels, while the adjacent plot illustrates the estimated intensity's standard error. The heatmap and 3D map further clarify the estimation, revealing the expected intensity inland from 2008 to 2013 and allowing analysis of variations across latitudes and longitudes.

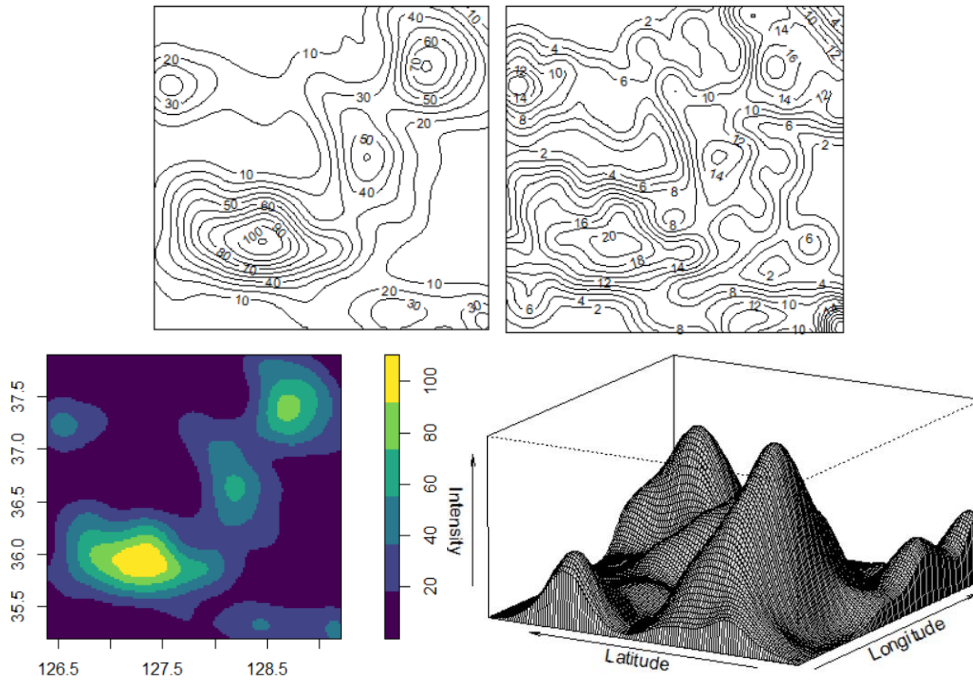


Figure 7. Contour plots, heatmap and 3D map of AGHs' occurrences intensity

2.3. Dependence of Intensity on Honeybees' Population

The Asian giant hornet predares honeybees intensively [38], so it is expected that the intensity of it is associated with the intensity of honeybees. In this study, *Apis mellifera* Linnaeus, 1758, a type of honeybees, are considered. They are highly populated in South Korea, and the occurrences records from 2008 to 2013 are shown in Figure 7. To investigate the reliance of *V. mandarinia*'s intensity on honeybees' population, the spatial distribution of bees' occurrences is applied and transformed to the density map as a spatial covariate. The spatial covariate is of interest in ecological studies because it may reflect the habitat preference of creatures [39]. The kernel-smoothing estimation method to attain the covariate was mentioned. It is necessary to embrace the same observation window so that the covariate map can adapt to the studied data of the hornet. Still, assuming the point pattern of honeybees in the window is consistent with the inhomogeneous Poisson point process, the intensity map shown in Figure 8 can be derived using Diggle's correction. Intuitively, there is a potential relationship between two insects' patterns in the window, and the statistical test called cumulative distribution function test [21] is used to test the dependency.

If one defines $C(x_i)$ as the covariate value at data point x_i , then

H_0 : The intensity does not depend on C ,

H_1 : The intensity depends on C in some unspecified mode.

Let $F_0(c) = \frac{|\{u \in W: C(u) \leq c\}|}{|W|}$ and $\hat{F}(c)$ be the cumulative distribution function of $C(x_i)$, and the discrepancy between two values is measured by Cramér-Von Mises statistic [40] $\omega^2 = n \int [\hat{F}(c) - F_0(c)]^2 dF_0(c)$, which turns out to be 5.6531, and the p-value of the test is effectively zero (smaller than 2.2×10^{-16} by computer), so the null hypothesis is rejected. Therefore, the dependence of *V. mandarinia*'s intensity on that of honeybees is significant statistically.

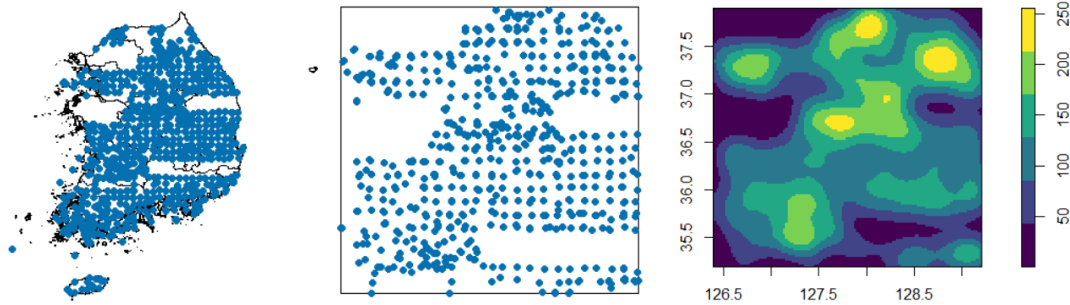


Figure 8. Cumulative occurrences map of honeybees in 2013 (left), occurrences pattern in observation window (Middle), intensity estimation plot (right)

Besides testing the dependence, the strength of the dependence should be explored because even though the dependence is statistically meaningful, the impact may be weak. The strength of the dependence can be measured by Receiver Operating Characteristic (ROC) and the area under the curve (AUC). In the ROC plot, the horizontal axis is $1 - F_0(c)$ which represents the proportion of the studied area satisfying $C(u) \geq c$, and the vertical axis is $1 - \hat{F}(c)$ which represents the fraction of data points with $C(x_i) \geq c$, or the fraction of data points of hornets falling in this area. The curve is shown in Figure 9, as c varies. The red dotted line serves as the ROC curve in the situation where the ratio is equivalent. Since the ROC curve lies substantially above the line, it is concluded that the level of dependence is positively strong to some extent. Reasonably, AUC can be interpreted as the probability that a randomly chosen *V. mandarinia*'s sighting location would have a greater value of honeybees' intensity than a randomly chosen location in the observation window, which turns out to be around 65.97%. In summary, the spatial intensity of Asian giant hornets has statistically convincing dependence on that of honeybees in South Korea of this research.

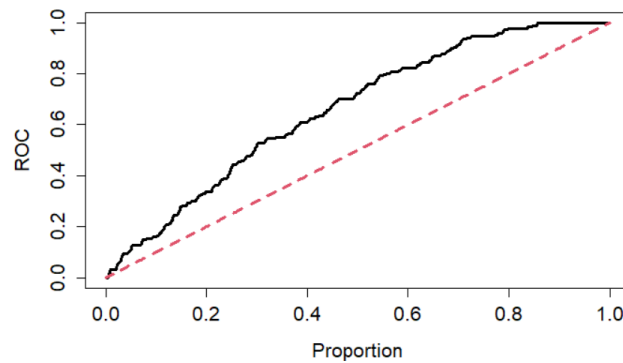


Figure 9. ROC curve ($AUC = 0.6597$)

3. Separated Spatial Pattern

3.1. Logic of Division

Unlike many creatures that are capable of gradually covering the living space with their population by reproduction annually, the hornet and bee's population do not spread over the continent. This is partly due to their colony cycle: take AGH as an example, the queen, after settling to nest and rear offspring [38], dies before winter [41]. During winter, new queens hibernate in the soil, creating seasonal gaps in occurrences, as verified in Figure 10. AGH's historical spatial distribution, as evidenced in Figures 11 and 12, is not expanding in the continent (Figures 1 and 2), though covering it overall. This subsection focuses on developing a robust model for the spatial distribution of AGH occurrences. Since sighting locations do not increase yearly, fitting the entire data set over years is impractical. Instead, the goal is to find a suitable point process model for annual data, contrasting with the non-parametric intensity modeling in the previous section.

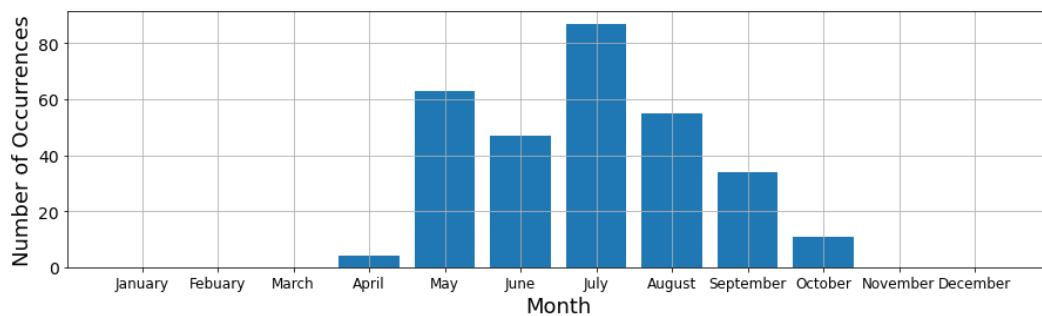


Figure 10. Accumulation of occurrences by month from 2008 to 2013

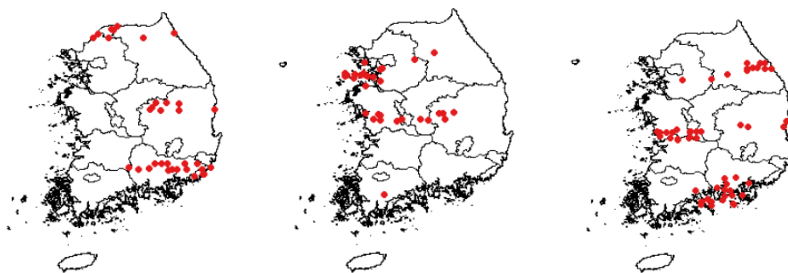


Figure 11. Annual maps of sighting records of *V. mandarinia* from 2008 to 2010 (left to right)

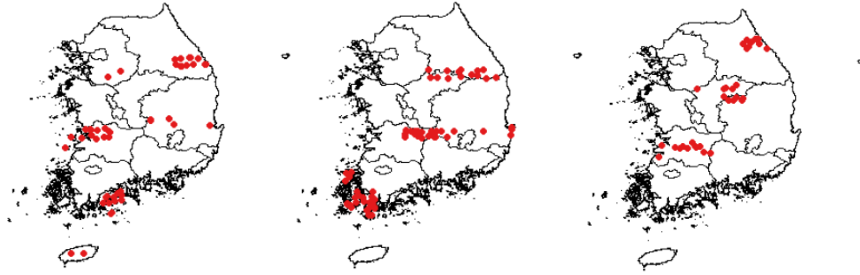


Figure 12. Annual maps of sighting records of *V. mandarinia* from 2011 to 2013 (left to right)

3.2. Methods and Ideas

In this paper, the Poisson point process is used in fitting the data for some reasons, and most importantly, the simplicity, since it is the basic model in spatial statistics. The queen of Asian giant hornets firstly finds a suitable place for nesting and breeding, and workers are going outside around the nest, making the occurrences scattering in some specific sites, so that the inhomogeneous Poisson process may be compatible with the pattern. For the homogeneous Poisson process, the intensity is not a function, but as the conclusion in the second section, here the pattern is assumed to be generated by the inhomogeneous Poisson process. To make the fitting reasonable, some assumptions are listed as followed for the AGH spatial pattern:

- The number of hornets sighting records in any region S in the window W follows a Poisson distribution.
- The expected number of records in any S in W is $\int_S \lambda(u) du$.
- For disjoint areas S_i in W , the number of records $n(X \cap S_i)$ are independent random variables.
- Let $n(X \cap S) = m$, then these m records are independent and identically distributed with probability density function $f(u) = \frac{\lambda(u)}{\int_S \lambda(x) dx}$

Another important tool for fitting here is the log-linear model [42], it has the general form

$$\lambda_\beta(u) = \exp(B(u) + \beta^T C(u)) = \exp(B(u) + \beta_1 C_1(u) + \beta_2 C_2(u) + \dots + \beta_m C_m(u))$$

where $B(u)$ is usually the log baseline, and $C(u)$ are covariates with parameters β .

For the parameter estimation, the maximum likelihood and numerical approximation are applied here. For a Poisson point process, the general likelihood function is given by

$$L(\lambda) = \frac{\lambda^{n(X)} \exp(-\lambda|W|)}{n(X)!},$$

where $n(X)$ denotes the number of points in the window, and $|W|$ the area of the window. Then for the inhomogeneous one it is given by

$$L(\beta) = \prod_{i=1}^{n(X)} \lambda_\beta(x_i) \cdot \frac{\exp(-\int_W \lambda_\beta(u) du)}{n(X)!},$$

where x_i is the observation points as mentioned before, and the loglikelihood is derived as

$$\log L(\beta) = \sum_{i=1}^{n(X)} \log(\lambda_\beta(x_i)) - \int_W \lambda_\beta(u) du.$$

However, the estimation is unable to be achieved by formal procedure, and here the By the quadrature numerical approximation specifically for inhomogeneous Poisson process [43], the final expression of the loglikelihood is given by

$$\log L(\beta) \approx \sum_{j=1}^m w_j \left[y_j \log(\lambda_\beta(u_j)) - \lambda_\beta(u_j) \right],$$

where $y_j = \frac{1}{w_j}$ for u_j as a data point, and 0 for u_j as a dummy point. The parameter β can be obtained by the method of fitting generalized linear models.

To build up the Poisson process model, the covariate used is the smoothing intensity estimation of honeybees in each year, and for fitting the data every year, the same observation window is fixed, which is the one used before. For instance, the spatial pattern of honeybees in 2012 is shown below in Figure 13. And it is hard to persuade one to believe that the pattern is formed by inhomogeneous Poisson point process or Cox process, but a linear work, so it may be invalid to search for the best σ by either mean square error or likelihood cross-validation methods and setting it as hyper-parameter is more appropriate. To deal with it, the optimized σ obtained in the second section in the article is used as an empirical constant for all the kernel-smoothing estimation of the intensity in each year, which is 0.1282. By applying this method, the covariate is drawn, and they are displayed in Figure 14 and 15.

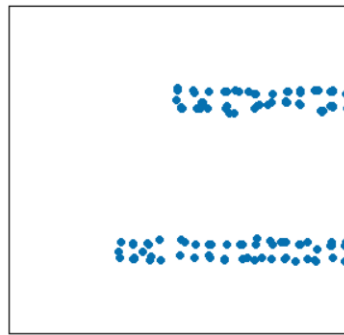


Figure 13. The point pattern of honeybees' occurrences in 2012

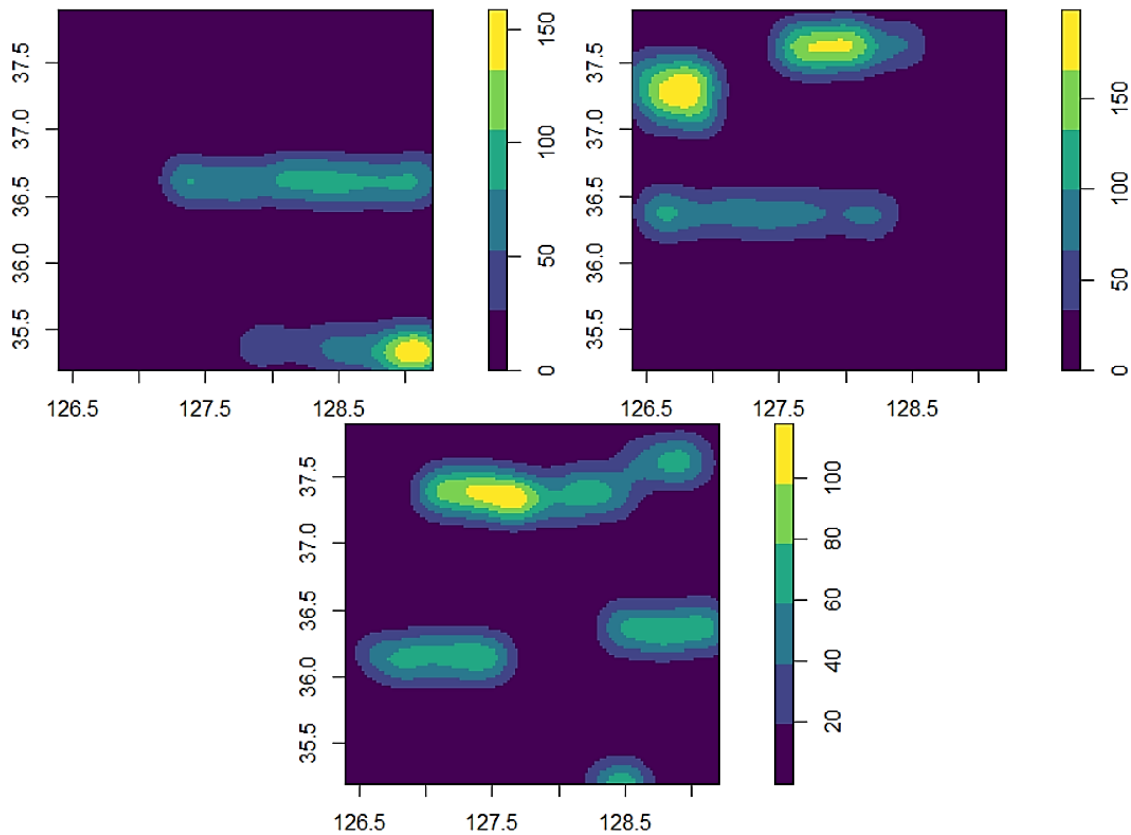


Figure 14. Occurrences pattern of honeybees from 2008 to 2010 (left to right)

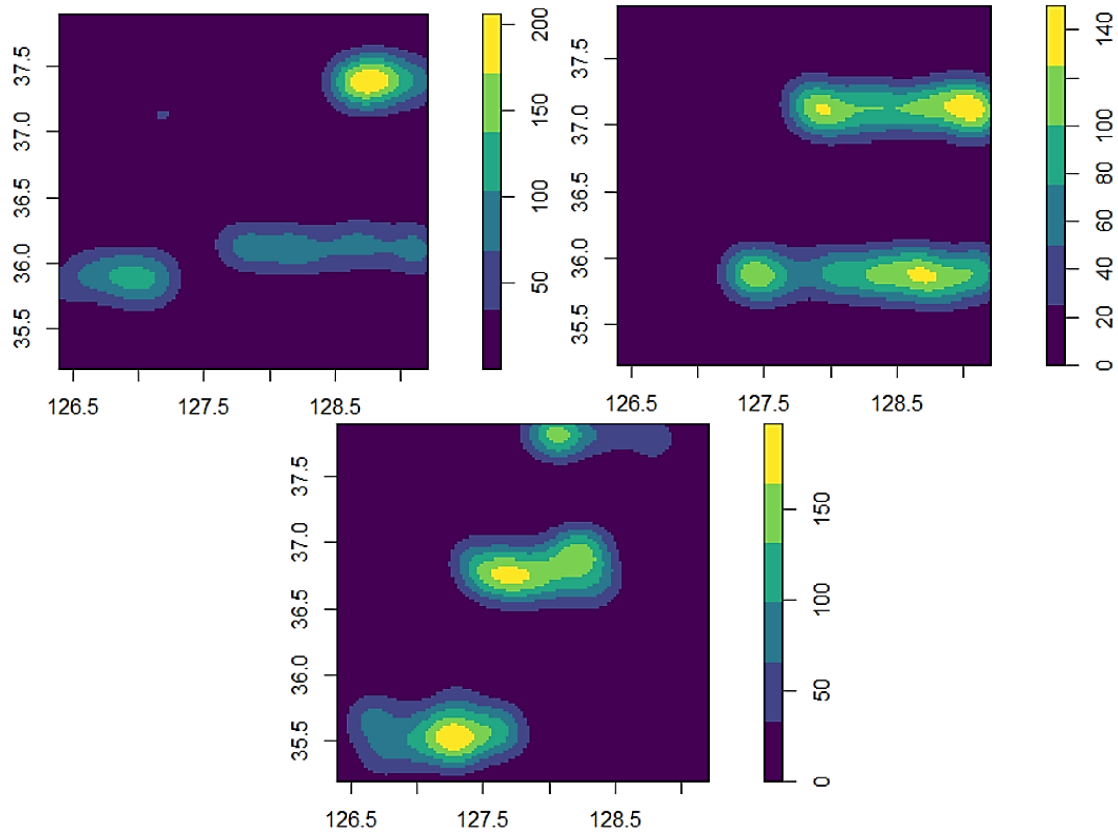


Figure 15. Occurrences pattern of honeybees from 2011 to 2013 (left to right)

3.3. Statistical Inference for Poisson Models

3.3.1. Model Parameter Test. In Section 2, the intensity of Asian giant hornets proves to depend on that of honeybees under the condition of non-parametric modelling, and it is rigorously necessary to test whether it still holds when establishing a Poisson point process model. The statistical test for this is called the likelihood ratio test (LRT) [44]. Suppose two models are built up, according to Section 3.2,

$$\lambda(u) = \exp(D);$$

$$\lambda(u) = \exp(D + \beta C(u)).$$

When fitting the first model, the parameter D means the average intensity $\bar{\lambda}$ per unit area, and when fitting the second model, where $C(u)$ is the value of honeybees' intensity derived by kernel-smoothing, with D as the intercept for $\log(\lambda(u))$. By giving the hypotheses $H_0: \beta = 0$ and $H_1: \beta \neq 0$. The likelihood ratio test statistic is given by

$$\Lambda = 2 \log \frac{L_2}{L_1}$$

where L_1 and L_2 are the maximum values of the likelihood function of the first and second model respectively, and the p-value is obtained by referring Λ to the χ^2 distribution with 1 degree of freedom. In order to prove the effect of the covariate in different numerical forms including $C(u)$, $C(u)^2$, $C(u)^{\frac{1}{3}}$, which would be used when modelling, the ones involved are tested with results shown in Table 1.

Table 1. Likelihood ratio test results

Year	$C(u)^2$	$C(u)$	$C(u)^{1/3}$
2008	3.423×10^{-13}	3.423×10^{-13}	$< 2.2 \times 10^{-16}$
2009	1.069×10^{-13}	1.069×10^{-13}	$< 2.2 \times 10^{-16}$
2010	8.602×10^{-16}	8.602×10^{-16}	$< 2.2 \times 10^{-16}$
2011	$< 2.2 \times 10^{-16}$	$< 2.2 \times 10^{-16}$	$< 2.2 \times 10^{-16}$
2012	$< 2.2 \times 10^{-16}$	$< 2.2 \times 10^{-16}$	$< 2.2 \times 10^{-16}$
2013	1.426×10^{-12}	1.426×10^{-12}	$< 2.2 \times 10^{-16}$

Please note that although the integral of the honeybees' population map is the positive, the density in some pixel may be under 0, which is the reason why $C(u)^{\frac{1}{3}}$ is chosen rather than $C(u)^{\frac{1}{2}}$. It can be found that all the tests are significant under great confidence, and the significance of $C(u)^{\frac{1}{3}}$ seems perform greatest for all the data. Under such circumstances, different numerical forms of the estimated intensity of honeybees would be put together in the log-linear model at the beginning for the model selection.

3.3.2. Model fitting and Selection. Before processing the model selection, some basic definitions are introduced. Wald test is the significance test of parameters in the log-linear model, and its test statistic is the same with that of traditional linear regression model, and this is because the parameter is calculated by the generalized linear models, as mentioned in 3.2. The Wald test statistic is given by $V = \frac{\hat{\beta}}{se(\hat{\beta})}$ where $\hat{\beta}$ is the maximum likelihood estimate. The other definition is the famous Akaike Information Criteria (AIC) [45] which is commonly used in applied statistics. It is also used in spatial statistics and is given by

$$AIC = -2\log(L(\beta)) + 2p,$$

where $L(\beta)$ is the maximum likelihood estimate, and p is the number of parameters in the model. People always prefers the model with lower AIC if it is the only standard chosen for model selection. Here, the backward stepwise subset selection method [46] is used for building optimized models. In this approach, a model containing all terms considered are analyzed initially, and then by deleting the variable one by one, the model with the smallest AIC score is achieved during the procedure, before confronting the situation where AIC cannot be reduced any more by withdrawing any other variable.

Firstly, the intuition of selecting variables should be stated. The coordinates of occurrences locations, i.e., the longitude and latitude of each record of AGH would be considered. As it can be visually found in Figure 16 and 17, the spatial distribution may have relevance with coordinates. For example, the points in the map of 2013 have gap vertically between clusters, and that may reflect on the latitude. An intercept would also be counted in models because it is traditional in running regression, and there seems to be no reason to drop it. The most important variable may be the kernel-smoothing intensity of honeybees, since it confirms to be significant in Poisson model in 3.3.1, and in the non-parametric modelling, the full spatial pattern of AGH relies heavily on it. Additionally, there would be adjustment of its scale. Here, $C(u)$, $C(u)^2$ and $C(u)^{\frac{1}{3}}$ will be considered together as mentioned, to strengthen the influence of the covariate. The article aims to find a model for each year's data.

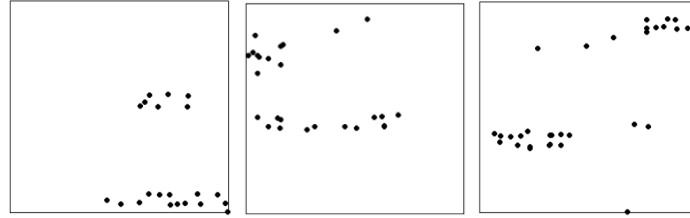


Figure 16. Spatial patterns simulation of V. mandarinia from 2008 to 2010 (left to right)

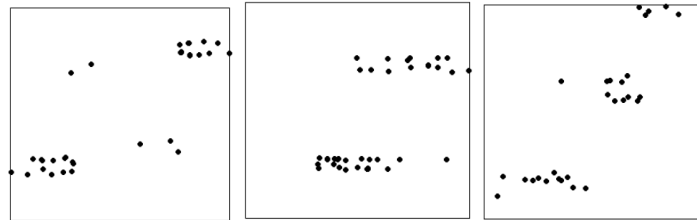


Figure 17. Spatial patterns simulation of V. mandarinia from 2011 to 2013 (left to right)

3.3.3. Report and Rethinking of the Model. The best fitting result is shown in Table 2, with estimate, standard error, and AIC to be shown. From the table, it can be concluded that for models fitted and selected by the smallest AIC score by each year's data, a great proportion of the parameters are statistically significant. Also, surprisingly it is found that the cubic root of the covariate is always included in the best model and is greatly significant, while the exact value of the covariate is included as well but with weaker significance, and the square of the covariate has never been included.

The task is to simulate the point pattern as good as possible by using a very simple and elementary model, so the explanation of some models may be accessible. For interpreting the parameter, it is easier for some models because they are explained in the form $\log(\lambda) = D + \beta_1 x_1 + \beta_2 x_2$. Take the model in 2008 as example, it can be estimated that for every one unit increase of the longitude in the window, the intensity of Asian giant hornets decreases about $|\exp(-1.1753) - 1| \cdot 100\% \approx 69\%$. Certainly, the longitude effect cannot outweigh the covariate effect, otherwise the fitting would be improper because the intensity function disables to generate points in the similar places like the factual plot.

Another problem which should be mentioned is the sign of the covariate effect to the response intensity in logarithm form according to the models fitted. Let

$$\log(\lambda(u)) = \alpha_1 C(u)^{\frac{1}{3}} + \alpha_2 C(u) + \alpha_3 C(u)^2 + \beta_1 \text{Latitude} + \beta_2 \text{Longitude} + \gamma.$$

If $\alpha_1 C(u)^{\frac{1}{3}} + \alpha_2 C(u) + \alpha_3 C(u)^2$ is defined as the overall impact of covariate to the logarithm of intensity, then the sign of the overall impact hardly becomes negative. For all models, firstly simplify the impact of covariate to the expression

$$f(x) = c_1 x^{\frac{1}{3}} - c_2 x,$$

where $c_1, c_2 > 0$ and $x \geq 0$. The derivative and second derivative are given by

$$\frac{df}{dx}(x) = \frac{1}{3} c_1 x^{-\frac{2}{3}} - c_2,$$

and

$$\frac{d^2 f}{dx^2}(x) = -\frac{2}{9} c_1 x^{-\frac{5}{3}},$$

respectively. Since $\frac{d^2f}{dx^2}(x) \leq 0$, the derivative is monotonically decreasing in $[0, +\infty)$. Therefore, $f(x)$ achieves its maximum at $x_0 = \left(\frac{3c_2}{c_1}\right)^{\frac{3}{2}} > 0$. $f(x)$ decreases monotonically when $x > x_0$, and there is a threshold $\left(\frac{c_1}{c_2}\right)^{\frac{3}{2}}$ when $f(x)$ becomes negative, which is obtained by setting $c_1 x^{\frac{1}{3}} = c_2 x$. Numerically, from 2008 to 2013, the thresholds are approximately 488.5711, 885.1234, 315.6818 and 982.4203, 530.1382, 511.1084 respectively, while even the maximum of the estimated intensity by Gaussian kernel-smoothing of the honey bees' complete point pattern is around 255.9695, which is smaller than these numbers, and it is expected that the real value cannot reach the theoretical threshold. Thus, the overall impact of the covariate being positive is verified in these models.

Table 2. Models report

Year	$C(u)^{\frac{1}{3}}$	$C(u)$	$C(u)^2$	Latitude	Longitude	Intercept	AIC
2008	4.5966* (0.1165)	-0.0741* (0.0372)	\	\	-1.1753** (0.3857)	31.6348* (13.5977)	-74.9495
2009	4.7845** (1.8269)	-0.0519 (0.0286)	\	\	-1.4431** (0.5056)	39.5744* (18.9962)	-97.5983
2010	15.9765** (5.0102)	-0.3446** (0.1151)	\	-0.3992 (0.2369)	\	12.6561 (32.4714)	-136.2609
2011	4.5064** (1.4338)	-0.0456* (0.0231)	\	-1.0950** (0.3755)	0.8031 (0.5260)	98.2772** (33.8322)	-204.8147
2012	14.3975** (5.2115)	-0.2198* (0.0884)	\	-1.2907*** (0.3056)	\	124.5854** (42.3676)	-223.3069
2013	6.7953** (2.1593)	-0.1063** (0.0367)	\	0.5461 (0.3674)	\	-87.2973 (48.1823)	-100.2346

3.3.4. Diagnostics of Models and Discussion. Diagnostics of models in each year are discussed here. In order to check whether the models' forms are correct, i.e. the assumption of inhomogeneous Poisson point process, the residual K-function [47] is constructed. Let x be a point pattern and suppose there is a summary function $F(x) = \sum_i f(x_i, x \setminus x_i)$ for any function f , note that $x \setminus x_i$ indicates the point pattern excluding x_i . For the empirical K-function mentioned in 2.2, it can be expressed by f in a general form

$$f(u, g) = \frac{1}{\bar{\lambda}(g \cup \{u\})^2 |W|} \sum_j e_K(u, g_j) I\{\|u - g_j\| \leq r\},$$

for any given $r \geq 0$, where u is a location and g is a point pattern. There is an equation called Campbell-Mecke formula [48], which states that

$$\mathbb{E}[\sum_{x \in X \cap S} f(x, X)] = \int_S \mathbb{E}[f(u, X)] \lambda(u) du,$$

where X is a Poisson point process with intensity $\lambda(u)$ and S is any region of it. Therefore, $\mathbb{E}[\sum_i f(x_i, x \setminus x_i)] = \int_W \mathbb{E}[f(u, x)] \lambda(u) du = \mathbb{E} \int_W f(u, x) \lambda(u) du$. If the residual function is defined as

$$RF(x) = \sum_i f(x_i, x \setminus x_i) - \int_W f(u, x) \lambda(u) du,$$

then swiftly the expectation of the residual function is given by

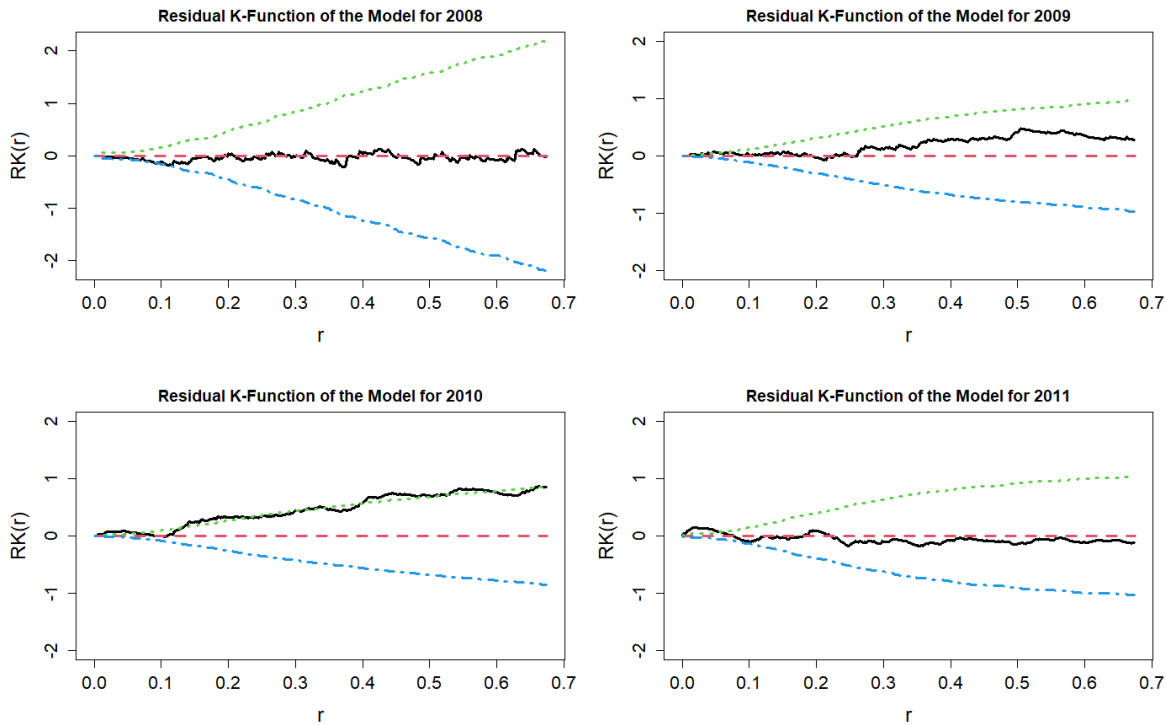
$$\mathbb{E}[RF(x)] = 0.$$

Since the definition of residual function is given above, the residual K-function can be defined as

$$RK_g(x) = \hat{K}_g(r) - \int_W f(u, g) \hat{\lambda}(u) du.$$

With this function, the principle is that if the model fitting is proper, $RK_g(x)$ is close to zero all the way, and the confidence interval of the residual is obtained by Poincare variance [47] that is simply of the form $\int_W f(u, g)^2 \hat{\lambda}(u) du$.

In Figure 18, the residual K-functions for the data of each year are plotted, the black solid line is the residual K-function obtained with the isotropic correction method, the red dotted line is the theoretical value, the space between the green and blue dotted line represents the 95% confidence interval of the function. One can find that in the figure, most plots of residual K-function are satisfactory since they are flowing close to zero, except the one for 2010, since it walks away from zero above the upper bound. Therefore, the model assumption for inhomogeneous Poisson point process is relatively successful for the data in the years 2008, 2009, 2011, 2012 and 2013, but seems to collapse for that in 2010, not too seriously. Returning to the plot in Figure 17, it seems that some points in 2010 are more likely to form a clustering shape than those of others when there are places for records. If the restriction is untied, the model for 2010 can still be useful, for the helpfulness to be shown in the lurking variable plot.



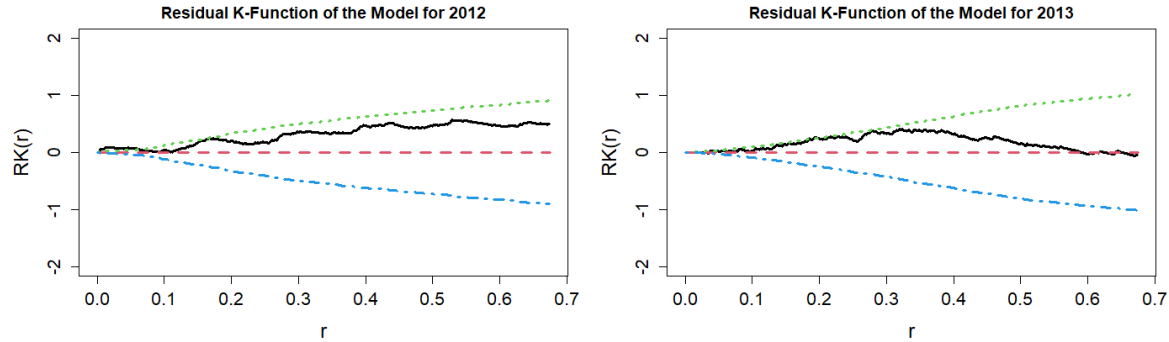
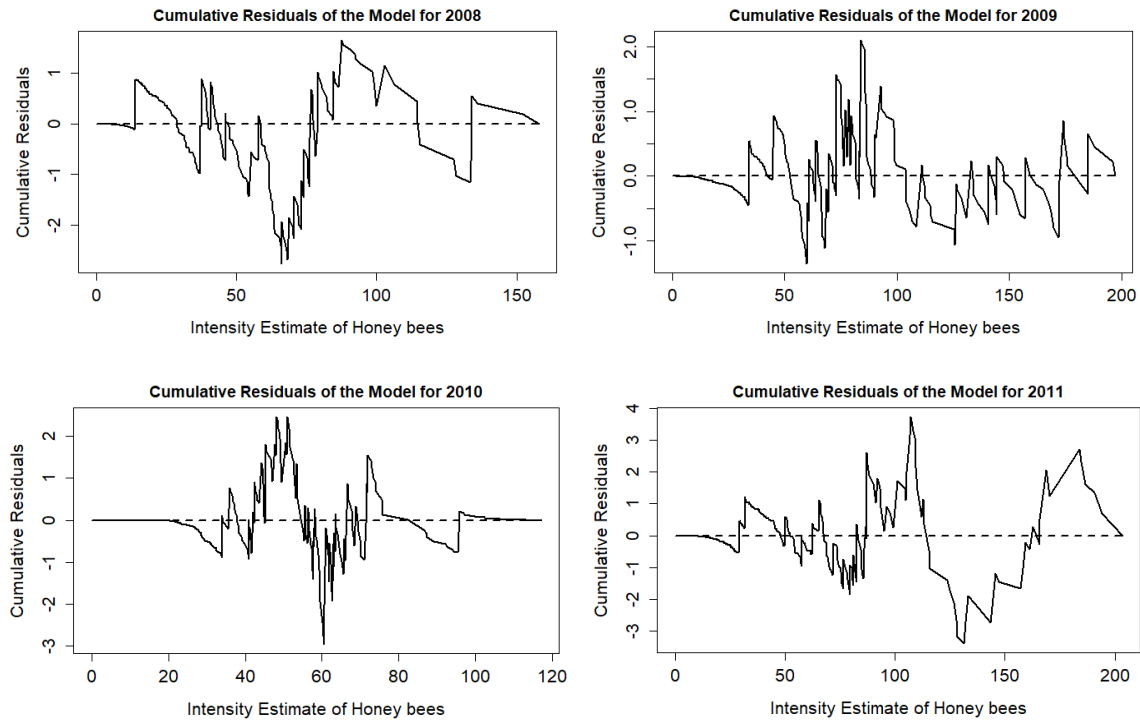


Figure 18. Plots of residual K-function of the models fitted for each year

Lurking variable plot [49], in this case, is a plot with horizontal axis on behalf of the value ρ , honeybees' intensity, and vertical axis performing the total residual for the space where the intensity of bees is no greater than ρ . Firstly, the point process residual [50] for the region S in X is defined as the difference between the number of observations and expected number derived by the fitted $\hat{\lambda}(u)$,

$$R(S) = n(X \cap S) - \int_S \hat{\lambda}(u) du.$$

Then the value of $R(Q(c)) = R(\{u \in W : C(u) \leq c\})$ is described as the lurking variable, where $C(\cdot)$ denotes the covariate, and W is the observation window. It can be interpreted as the difference between the cumulative number of Asian giant hornets sighted on locations where bees' intensity is less or equal to some value, and the expected number according to the model. From Figure 19, it can be noticed that most of the absolute values of cumulative residuals rarely go beyond around 4, given that the average sighted number per year is around 50, so the models fitted are relatively reliable.



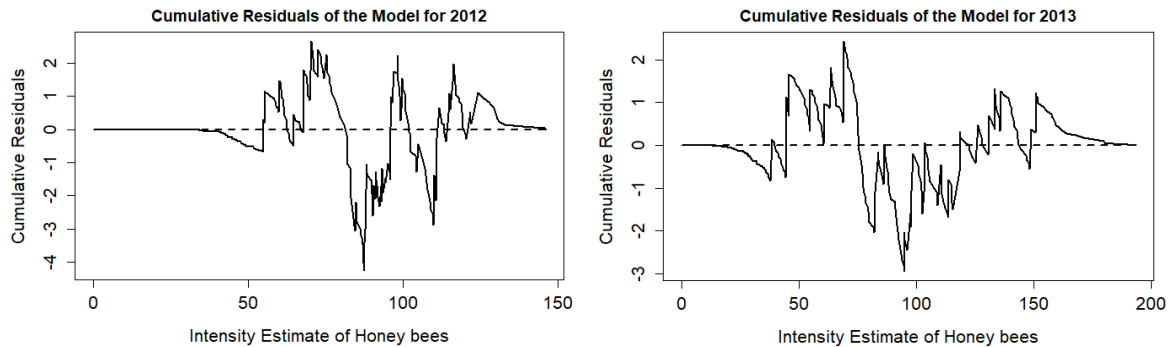


Figure 19. Plots of cumulative residuals of models fitted for each year

4. Conclusion

From the data sampled from 2008 to 2013, the intensity of Asian giant hornets, *Vespa Mandarinina Smith, 1852* proves to be highly dependent on the intensity of the honeybee which is *Apis mellifera Linnaeus, 1758* in the study. By splitting the spatial pattern of Asian giant hornets by year, the intensity of the hornet can be modelled by inhomogeneous Poisson point process with log-linear model, setting the intensity of honeybees as covariate that is estimated by Gaussian kernel-smoothing method. Almost all the models established are convincing, with parameters to be statistically significant, and the entire group of the point process is appropriate and valid for simulation. The expectation of inferring the spatial pattern of Asian giant hornets by a rudimentary statistical model has been achieved in the paper.

References

- [1] Bérubé, C 2020 Giant Alien Insect Invasion Averted – Canadian Beekeepers Thwart Apicultural Disaster *American Bee Journal* 160(2) pp 209–214
- [2] Alaniz, A J, Carvajal, M A and Vergara, P M 2021 Giants are coming? Predicting the potential spread and impacts of the giant Asian hornet (*Vespa mandarinia*, Hymenoptera: Vespidae) in the USA *Pest Manag Sci* 77(1) pp 104–112
- [3] Hirano, K, Tanikawa, A 2020 Ocular injury caused by the sprayed venom of the Asian Giant Hornet (*Vespa mandarinia*) *Case Reports in Ophthalmology* 11(2) pp 430–435
- [4] Park, M, Zhang D and Landau, E 2013 Deadly giant hornets kill 42 people in China CNN October 2013 Available at: <https://edition.cnn.com/2013/10/03/world/asia/hornet-attack-china/index.html>
- [5] Yamane, S 1976 Morphological and taxonomic studies on vespine larvae, with reference to the phylogeny of the subfamily Vespinae (Hymenoptera: Vespidae), *Insecta Matsumurana, Entomology* 8 pp 1–45
- [6] Osterloff, E 2018 Why Asian hornets are bad news for British bees CNN, October 2013. Available at: <https://www.nhm.ac.uk/discover/why-asian-hornets-are-bad-news-for-british-bees.html>.
- [7] Zhu, G, Gutierrez, J I, Looney, C, Crowder, D W 2020 Assessing the ecological niche and invasion potential of the Asian giant hornet *Proceedings of the National Academy of Sciences* 117 (40) pp 24646–24648
- [8] Diggle, P J 2013 *Statistical Analysis of Spatial and Spatio-Temporal Point Patterns*. 3rd ed (London: Chapman and Hall/CRC) chapter 4 p 55
- [9] Kerscher, M 2000 Statistical analysis of large-scale structure in the Universe, in Mecke, K.R. and Stoyan, D. (eds.) *Statistical Physics and Spatial Statistics, Lecture Notes in Physics* Berlin: Springer-Verlag pp 36–71
- [10] Law, R, Illian, J, Burslem, D F R P, Gratzer, G, Gunatilleke, C V S, Gunatilleke, I A U N (2009) Ecological information from spatial patterns of plants: insights from point process theory *Journal of Ecology* 97(4) pp 616–628

- [11] Glasbey, C, and Roberts, I M 1997 Statistical analysis of the distribution of gold particles over antigen sites after immunogold labelling *Journal of Microscopy* 186(3) pp 258–262
- [12] Pfeiffer, D U, Robinson, T P, Stevenson M, Stevens K B, Rogers D J, Clements A C A 2008 *Spatial Analysis in Epidemiology* (Oxford: Oxford University Press)
- [13] Penttinen, A 1984 *Modelling Interaction in Spatial Point Patterns: Parameter Estimation by the Maximum Likelihood Method* vol 7 (Jyväskylä: Jyväskylän yliopisto)
- [14] Chainey, S and Ratcliffe, J 2005 *GIS and Crime Mapping* (Oxfordshire: Taylor & Francis) pp 79–110
- [15] Kim, J K, Choi, M B and Moon, T Y 2006 Occurrence of *Vespa velutina* Lepeletier from Korea, and a revised key for Korean *Vespa* species (Hymenoptera: Vespidae *Entomological Research* 36(2), pp 112–115
- [16] Martin, S J 1990 Nest thermoregulation in *Vespa simillima*, *V. tropica* and *V. analis*, *Ecological Entomology* 1990(15) pp 301–310
- [17] Kim, J K and Kim, I K 2011 Discovery of *Vespa binghami* (Vespidae: Hymenoptera) in Korea, *Animal Systematics Evolution and Diversity* 27(1) pp 105–107
- [18] Moo-Llanes, D A 2021 Inferring Distributional Shifts of Asian Giant Hornet *Vespa mandarinia* Smith in Climate Change Scenarios *Neotropical Entomology* 50(4) pp 673–676
- [19] Nuñez-Penichet, C., Osorio-Olvera, L, Gonzalez, V, Cobos, M E , Jiménez, J L, DeRaad, D, Alkische, A, Contreras-Díaz, R, Nava B A, Utsumi, K, Ashraf, U, Adeboje, A., Peterson, A, Soberón, J 2021 Geographic potential of the world's largest hornet, *Vespa mandarinia* Smith (Hymenoptera: Vespidae), worldwide and particularly in North America *PeerJ* 9(1) e10690
- [20] Elith, J, Leathwick, J R (2009) Species Distribution Models: Ecological Explanation and Prediction Across Space and Time *Annual Review of Ecology Evolution, and Systematics* 40 (1) pp 677–697
- [21] Berman, M 1986 Testing for spatial association between a point process and another stochastic process *Journal of the Royal Statistical Society, Series C (Applied Statistics)* 35(1) pp 54–62
- [22] Funwi-Gabga, N and Mateu, J 2012 Understanding the nesting spatial behaviour of gorillas in the Kagwene Sanctuary, Cameroon *Stochastic Environmental Research and Risk Assessment* 26(6) pp 793–811
- [23] Kwon, Y National Institute of Ecology 2022 'niek_2022' Available at: https://www.gbif.org/occurrence/search?country=KR&dataset_key=1e09d985-8f58-4a34-ab86-17e55f203c8a&taxon_key=5871429&year=2008,2013 (Accessed 26 November 2023)
- [24] Kwon, O and Choi, M B 2020 Interspecific hierarchies from aggressiveness and body size among the invasive alien hornet, *Vespa velutina nigrithorax*, and five native hornets in South Korea *PLoS ONE* 15(7) e0226934
- [25] Gelfand, A E, Diggle, P J, Fuentes, M, Guttorp, P 2010 *Handbook of Spatial Statistics* (London: Chapman and Hall/CRC) p 354
- [26] Kojima, J National Museum of Nature and Science, Japan 2021 Insect Specimens of Ibaraki University (Faculty of Science). Version 1.4 Available at: https://www.gbif.org/occurrence/search?dataset_key=04bcb07-811a-4870-89ec-4055fa88ea8d&taxon_key=5871429&year=1000,2022 (Accessed 26 November 2023)
- [27] Baddeley, A, Rubak, E and Turner, R 2015 *Spatial point patterns methodology and applications with r* (Florida: Chapman and Hall/CRC) pp 19–48 173
- [28] Illian, J, Penttinen, A, Stoyan, H, Stoyan, D 2008 *Statistical Analysis and Modelling of Spatial Point Patterns* (Chichester: John Wiley & Sons Ltd) pp 4–5
- [29] Strauss, D J 1975 A model for clustering, *Biometrika* 62(2) pp 467–475
- [30] Dahlhaus, R and KÜNSCH, H 1987 Edge effects and efficient parameter estimation for stationary random fields *Biometrika* 74(4) pp 877–882
- [31] Ripley, B D 1977 Modelling spatial patterns (with discussion) *Journal of the Royal Statistical Society, Series B* 39(2) pp 172–212

- [32] Horvitz, D G and Thompson, D J 1952 A generalization of sampling without replacement from a finite universe *Journal of the American Statistical Association* 47(260) pp 663–685
- [33] Bithell, J F 1990 An application of density estimation to geographical epidemiology *Statistics in Medicine* 9(6) pp 691–701
- [34] Diggle, P J 1985 A kernel method for smoothing point process data *Journal of the Royal Statistical Society, Series C (Applied Statistics)* 34(2) pp 138–147
- [35] Berman, M and Diggle, P J 1989 Estimating weighted integrals of the second-order intensity of a spatial point process *Journal of the Royal Statistical Society, Series B* 51(1) pp 81–92
- [36] Loader, C 1999 *Local Regression and Likelihood* (New York: Springer) pp 90–92
- [37] Daley, D J and Vere-Jones, D 1988 *An Introduction to the Theory of Point Processes* (New York: Springer-Verlag) vol 1
- [38] Matsuura, M and Sakagami, S F 1973 A Bionomic Sketch of the Giant Hornet, *Vespa mandarinia*, a Serious Pest for Japanese Apiculture (With 12 Text-figures and 5 Tables) *Journal of the Faculty of Science Hokkaido University Series VI. Zoology (Hokkaido University)* 19(1) pp 125–162
- [39] Manly, B J F, McDonald, L L and Thomas, D L 1993 *Resource Selection by Animals: Statistical Design and Analysis for Field Studies* (London: Chapman and Hall/CRC) pp 196–197
- [40] Cramér, H 1928 On the composition of elementary errors: II, Statistical applications' *Scandinavian Actuarial Journal* 11(1) pp 141–180
- [41] Matsuura, M 1984 Comparative biology of the five Japanese species of the genus *Vespa* (Hymenoptera, Vespidae) *The Bulletin of the Faculty of Agriculture - MIE University* 69 pp 1–131
- [42] Cox, D R 1972 The statistical analysis of dependencies in point processes *Stochastic Point Processes: Statistical Analysis, Theory, and Applications* (New York: Wiley) pp 55–66
- [43] Berman, M and Turner, T R 1992 Approximating point process likelihoods with GLIM *Journal of the Royal Statistical Society, Series C (Applied Statistics)* 41(1) pp 31–38
- [44] Rathbun, S L and Cressie, N 1994 Asymptotic properties of estimators of the parameters of spatial inhomogeneous Poisson point processes *Advances in Applied Probability* 26(01) pp 122–154
- [45] Akaike, H 1974 A new look at the statistical model identification *IEEE Transactions on Automatic Control* 19(6) pp 716–723
- [46] Derksen, S and Keselman, H J 1992 Backward, forward and stepwise automated subset selection algorithms: frequency of obtaining authentic and noise variables *British Journal of Mathematical and Statistical Psychology* 45(2) pp 265–282
- [47] Baddeley, A, Rubak, E and Møller, J 2011 Score, pseudo-score and residual diagnostics for spatial point process models *Statistical Science* 26(4) pp 613–646
- [48] Mecke, J 1967 Stationäre zufällige maße auf lokalkompakten abelschen Gruppen *Probability Theory and Related Fields* 9(1) pp 36–58
- [49] Baddeley, A, Turner, R, Møller, J, Hazelton, M 2005 Residual analysis for spatial point processes (with discussion) *Journal of the Royal Statistical Society, Series B* 67(5) pp 617–666
- [50] Baddeley, A, Møller, J and Pakes, A G 2008 Properties of residuals for spatial point processes, *Annals of the Institute of Statistical Mathematics* 60(3) pp 627–649