

MT4537 - Spatial Statistics

Project 1 - Simulation and Model Fitting

Contents

1	Simulation	3
1.1	Description of the Simulation Process	3
1.1.1	Usability and Reproducibility	3
1.1.2	Parent Points	4
1.1.3	Thinning	5
1.2	Realisations from the Simulation	7
1.3	Fitting Thomas and Matern Process Models	8
1.3.1	Minimum Contrast Method	8
2	Model Fitting	10
2.1	Kernel Smooths	10
2.2	Fitting Using Covariates	12
2.2.1	AIC	12
2.2.2	Using Coordinates	12
2.2.3	Using Altitude and Gradient	12
A	R Code	14
A.1	Algorithm Example	14
A.2	Simulation Functions	15

Chapter One

Simulation

1.1 Description of the Simulation Process

This section provides a description of the process followed to create the simulation algorithm for a Thomas cluster process by deriving a random driving intensity and then using thinning to create a given realisation. For reference, the code for the simulation algorithm is contained in the *src/simul_funcs.R* file.

1.1.1 Usability and Reproducibility

The simulation algorithm is written using functions within *R* since this allows the simulations to be re-run multiple times easily and without having to repeat code. Making use of functions allowed inputs to be created that let the user control some elements of the simulations. These inputs included the mean of the distribution of the number of parents, the bandwidth used in creating the driving intensity and a random seed.

Initially, there is a set of input checks which are just there to improve usability of the function. It ensures the user inputs are valid to reduce the chance of unexplained errors later in the process. The data type and limitations for each input is described above the function. The next step was to set the random seed. This allows the results of the simulation (ie. the realizations) to be reproduced if the same inputs are given to the function. If this was not included then the realizations would be different every time the function is run, even if all the other inputs remained the same.

1.1.2 Parent Points

Moving on to the distribution of the "parent" points, the theory of Neyman-Scott processes assumes that the parent points are generated from a homogeneous Poisson process. This means the number of these parent points can be given by a single random draw from a Poisson distribution. This is exactly what has been implemented in the simulation algorithm where the mean of this distribution is controlled by a user input which can be manipulated as required.

The next step is to place the required number of parent points randomly and independently across the region. Since these parents come from a homogeneous process they have equal probability of being placed anywhere within the region which implies that they can be placed using random draws from a uniform distribution. This is done in the algorithm by taking the given number of random draws from a uniform distribution with limits given by the length of the region to give the x -coordinates and another set of similar random draws giving the y -coordinates. These random draws were stored as vectors with equal length and converted to a *spatstat* *ppp* object for use in the next step.

A driving intensity was then derived using the randomly located parent points as the centres of Gaussian densities which were symmetric in every direction. These densities have a single parameter σ which is the standard deviation of the normal distribution being placed at each parent point. This can be controlled by the user and changed to fit the situation being simulated. This driving intensity was created using the *density* function from *spatstat* using Gaussian kernels and the *sigma* parameter as described above.

Now we will consider a short example of this part of the algorithm. The *R* code used for this example is contained in Section A.1 of the Appendix. First, set the seed to 123 and simulate a value from a Poisson distribution with mean 11. This returns a value of 9 parent points. The next step is to generate the (x, y) -coordinates for each of these 9 points using random draws from a uniform distribution in the range $[0, 1]$. This returns vectors for the x -coordinate and the y -coordinate, the values of which are contained in Table 1.1. These points can be plotted as seen on the left in Figure 1.1. The final step in this section of the algorithm is to devise the driving intensity from these parent points by placing symmetric Gaussian densities at these parent points with a standard deviation of 0.1. This intensity can be seen on the right side of Figure 1.1.

x	0.409	0.883	0.940	0.046	0.528	0.892	0.551	0.457	0.957
y	0.453	0.678	0.573	0.103	0.900	0.246	0.042	0.328	0.955

Table 1.1: x - and y -coordinate values used in the example in Section 1.1.2.

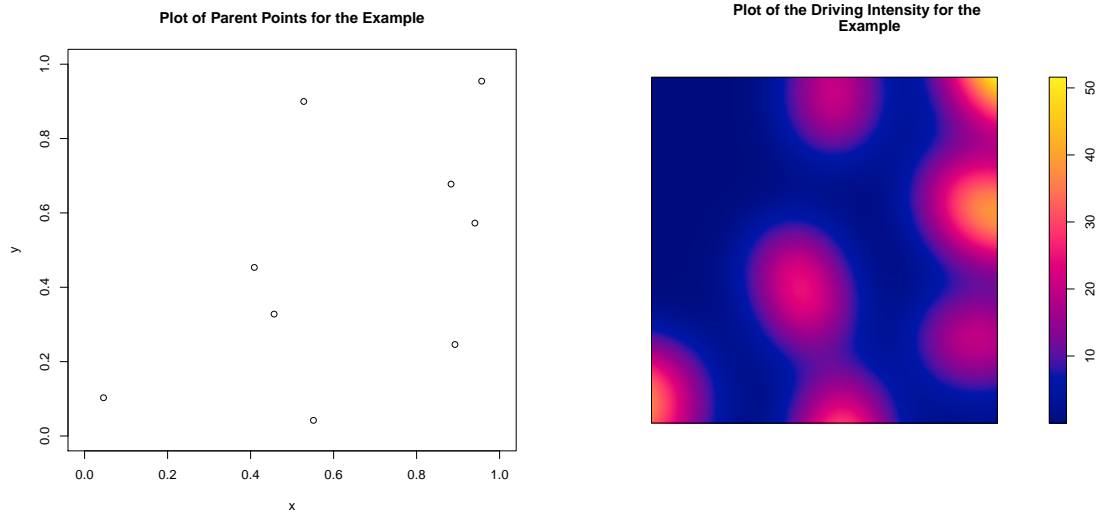


Figure 1.1: This contains the plots for the example in Section 1.1.2. The left hand plot is the parent points and the right hand plot is the resulting driving intensity using Gaussian densities and a standard deviation of 0.1. This was completed in a unit square.

1.1.3 Thinning

The thinning section of the algorithm took part in a separate function to avoid the main simulation function becoming too crowded. This function took as inputs the driving intensity based on the parent points derived in the previous step and the length of the sides of the square being simulated in.

This thinning process is completed by first generating a sample from a homogeneous Poisson process using the *rpoispp()* function from the *spatstat* library this time. The intensity for this homogeneous Poisson process is given by the maximum intensity value observed in the derived driving intensity that came from the parent points and is generated in the same window as specified earlier.

The thinning step is then applied to this realisation of the homogeneous Poisson process. Two empty vectors are created to store the (x, y) -coordinates that are to kept as points in the final simulation. The (x, y) -coordinates are extracted from the *ppp* object created during the simulation from the homogeneous Poisson. All the values in the previously derived driving intensity are transformed to probabilities by dividing each by the maximum value of the intensity observed so they are all now in the range $[0, 1]$. Now, for each set of (x, y) -coordinates, a Bernoulli random variable is simulated. This has the following PDF:

$$f(x) = \begin{cases} 1 & \text{with probability } p \\ 0 & \text{with probability } 1 - p \end{cases}$$

The probability p of returning a value of 1 is given by the value of the transformed density evaluated at the point under consideration. The point was kept and added to the initially empty storage vectors if a value of 1 is returned by the Bernoulli trial and is removed (not added to the storage vectors) from the simulation if 0 is returned.

The (x, y) -coordinates that are kept after the thinning process described above are then converted to a *ppp* object then returned. The overall simulation function then returns a list containing these simulation points as a *ppp* object and the driving intensity derived in Section 1.1.2.

Continuing the example from the previous section, the full simulation from a homogeneous Poisson process is created. This is shown on the left of Figure 1.2. The thinning process is then applied making use of the driving intensity depicted on the right of Figure 1.1. The plot on the right of Figure 1.2 shows all the points again but those points in red are the points that are kept as part of the simulation. All the other points are removed from the simulation.

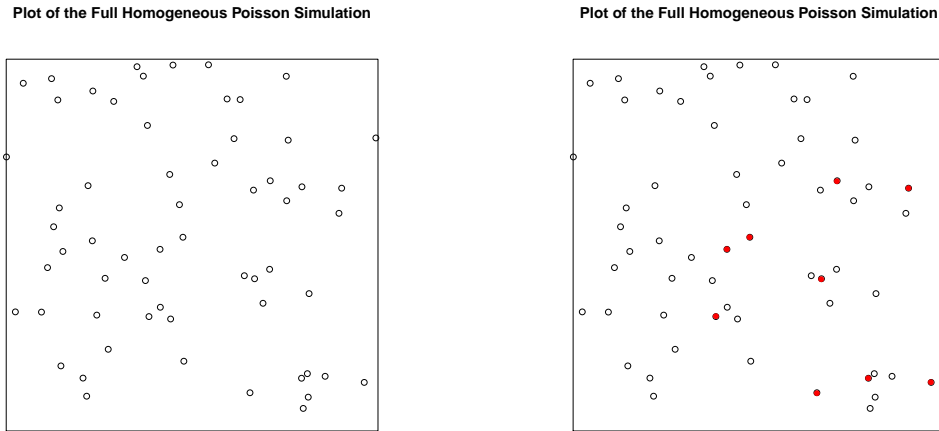


Figure 1.2: This contains the plots for the example in Section 1.1.3. The left hand plot is the simulation from the homogeneous Poisson process and the right hand plot is the thinned simulation where the red points were those that were kept as points in the simulation (the rest were removed). This was completed in a unit square.

1.2 Realisations from the Simulation

This section contains two realisations from the simulation algorithm described in the previous section. Figure 1.3 shows the plots from the first simulation where the mean number of parent points was 30, the standard deviation of the Gaussian densities placed at the parent points was 0.1, a random seed of 200 and the simulation was completed in the unit square. The plots for the second simulation are contained in Figure 1.4. This simulation was completed using a mean number of parents of 45, standard deviation of the Gaussian densities equal to 0.05 and a random seed of 465. This was again simulated within the unit square.

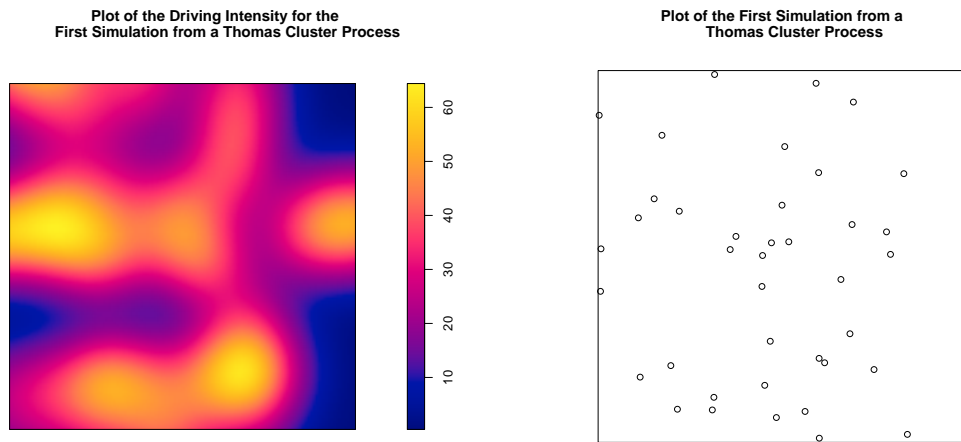


Figure 1.3: This contains the plots for the first realisation from the simulation of a Thomas cluster process. The left hand plot is the driving intensity and the right hand plot is the resulting set of spatially distributed points. This was simulated in a unit square.

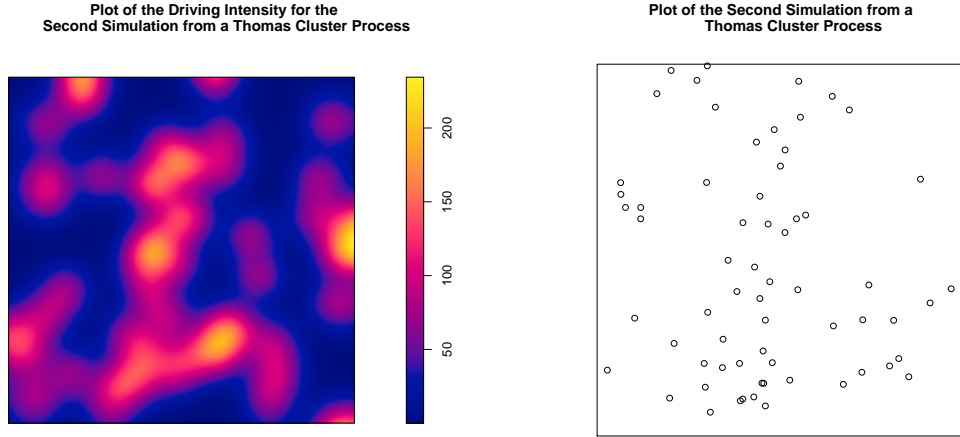


Figure 1.4: This contains the plots for the second realisation from the simulation of a Thomas cluster process. The left hand plot is the driving intensity and the right hand plot is the resulting set of spatially distributed points. This was simulated in a unit square.

1.3 Fitting Thomas and Matern Process Models

1.3.1 Minimum Contrast Method

This is a common method used for estimating parameters for Poisson cluster processes which is based on the idea of least-squares estimation. It makes use of an estimated and theoretical summary characteristic and tries to find the parameter values which minimise the difference between these. The method can be seen as trying to minimise

$$\Delta(\theta) = \int_{s_1}^{s_2} |\hat{S}(r) - S_{\theta}(r)|^{\beta} dr$$

with respect to θ where this is a vector of the parameters. The vector of parameters that minimise this expression are the estimators of the parameters. $S(r)$ is a summary characteristic depending on the distance from a point r . The other values included in the expression are generally arbitrary. However, it is normal practice to select $\beta = 2$. The limits of the integration expression remain arbitrary with no standard selection for these.

In reality, the formula given above for the minimum contrast method is discretised (broken up into discrete points) to make a summation formula. This means the method now tries to minimise using

$$\int_{s_1}^{s_2} |\hat{S}(r) - S_{\theta}(r)|^{\beta} dr \approx \sum_{i=0}^k |\hat{S}(\gamma_i) - S_{\theta}(\gamma_i)|^{\beta}$$

where $S(r)$ and β remained defined as above. Meanwhile, $\gamma_0 = s_1$, $\gamma_k = s_2$ and $\gamma_i = \gamma_0 + i\delta$ where $\delta = \frac{s_2 - s_1}{k}$ and k is a positive integer. This is just breaking the range of values for r into $k + 1$ discrete points to evaluate the summary characteristic function $S(r)$ at to give the total difference. Taking the limit as k tends to ∞ of the sum will lead back to the initial integration definition for the method.

The summary characteristics commonly used are either

$$S(r) = K(r) \quad \text{or} \quad S(r) = g(r)$$

where $K(r)$ is Ripley's K-function and $g(r)$ is the pair correlation function. Throughout this section the minimum contrast method will be used with $g(r)$. This is defined in terms of $K(r)$ and contains the same information as it. $K(r)$ is defined to be the mean number of distinct points found within a distance r from any given point (not including the point itself). $g(r)$ is then a summary characteristic defined to be proportional to the derivative of $K(r)$ defined previously. This can be written out as

$$g(r) = \frac{K'(r)}{2\pi r} \quad \text{for } r \geq 0$$

Chapter Two

Model Fitting

The *bei* data set contains information on the location of 3604 trees and *bei.extra* provides additional covariates showing the altitude across the region and also the gradient. A plot of the *bei* data is provided on the left of Figure 2.1. Here we are trying to fit a model to describe the location of the trees across the region.

2.1 Kernel Smooths

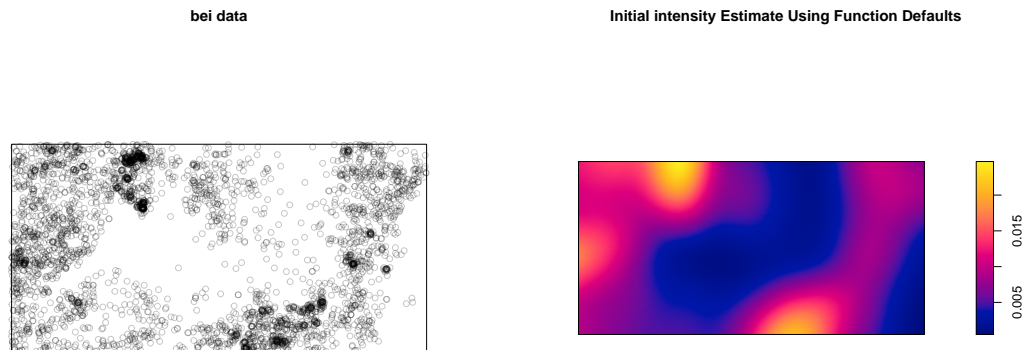


Figure 2.1: The plot on the left is the *bei* data set. On the right is an initial estimate of the intensity for this data using the *density* function with the arguments left at the default values.

The first step when trying to model this data is to fit a kernel smooth. This is fitted using the *density.ppp* function provided in *spatstat*. This function is simply placing symmetric kernels at every point in the pattern and then summing all the values of these kernels at every point. The shape and spread of the kernel can be controlled by specifying different types of kernels and specifying the bandwidth (*sigma*).

The initial smooth, depicted on the right of Figure 2.1, is fitted using the function defaults. It is difficult to tell with this many points which kernel is being used. Since it appears to be relatively smooth, it is unlikely to be the disc kernel but it is not easy to decide between the other three options (these being Gaussian, quartic and Epanechnikov). Looking at the help files for the function, the default kernel used is in fact the Gaussian. Again, just looking at the plot it is difficult to make any kind of guess as to the value being used for the bandwidth. However, the help files says that there is a "simple rule of thumb that depends only on the size of the window" to calculate a value for *sigma* when none is specified by the user. In this case, it appears to have a value of 50 and a window size of 1000 by 500 (this was concluded using trial and error). This implies the formula being used is either $0.05a$ or $0.1b$ where a is longer than b but this is only based on this one example.

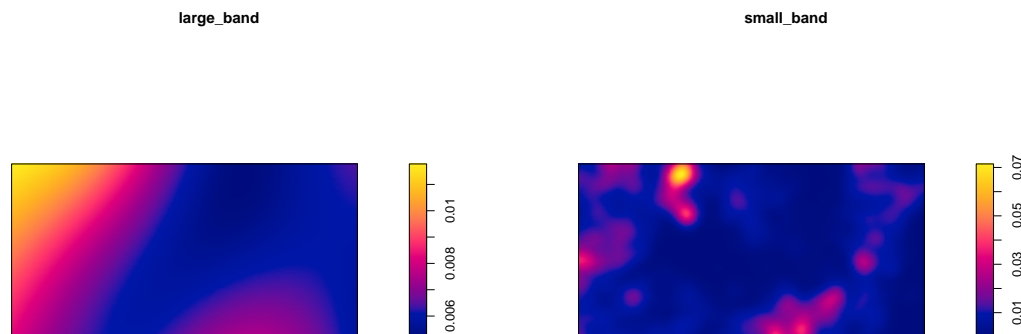


Figure 2.2: Initial estimate of the intensity for the *bei* data using the *density* function with Gaussian kernels but σ equal to 200 on the left and 20 on the right.

Figure 2.2 provides two markedly different intensity estimates. These are both also markedly different from the initial estimate using the function defaults. The only change in these intensities is the use of different values for the bandwidth. It is clear to see that

the higher bandwidth leads to a much flatter intensity while the smaller value leads to a much more bumpy estimate of the intensity.

2.2 Fitting Using Covariates

This section covers the model fitting process used when adding covariates into the models through the *ppm* function in *spatstat*. The models being fitted here are nearly equivalent to a Poisson GLM. They use a log-link since the intensity being estimated must be strictly positive.

2.2.1 AIC

The model selection criteria used was Akaike's Information Criterion (AIC) which is calculated as follows:

$$AIC = -2\log[L(\hat{\theta})] + 2q$$

where $L(\hat{\theta})$ is the likelihood function evaluated at its maximum which are the outputs of the modelling functions. q is defined to be the number of parameters estimated in the model and acts as a penalty for estimating more parameters. Model selection using AIC is carried out by selecting the model which has the lowest AIC value.

2.2.2 Using Coordinates

The first model fitted was just a constant model using the formula " $bei \sim 1$ ". Looking at the intensity estimated in Section 2.1 this will not be a good model which should be evident in the AIC values but this model is included for completeness. The next model fitted was one that depended on x and y individually using the formula " $bei \sim x + y$ " which was followed by the model based on the interaction term using the formula " $bei \sim x * y$ ". The final three models fitted were all polynomials of degree 2, 3 and 6 respectively. These polynomial models include coefficients for all the powers of the coordinates up to and include the degree given. For example, the form of the polynomial model of degree 3 is

$$\log(\lambda(x, y)) = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3 + \beta_4y + \beta_5y^2 + \beta_6y^3$$

where the β_i s are the coefficients to be estimated in the model.

The AIC values for all these models are contained in Table 2.1

2.2.3 Using Altitude and Gradient

Model Formula	AIC
1	42763.92
$x + y$	42545.68
$x * y$	42455.77
$\text{poly}(x, 2) + \text{poly}(y, 2)$	42234.23
$\text{poly}(x, 3) + \text{poly}(y, 3)$	42157.51
$\text{poly}(x, 6) + \text{poly}(y, 6)$	41966.5

Table 2.1: AIC values for the models using the (x, y) -coordinates.

Appendix A

R Code

A.1 Algorithm Example

Code used for the example in Section 1.1.2 and 1.1.3

```
library(spatstat)
```

```
set.seed(123)
```

```
num_p <- rpois(1, 11)
```

```
x_p <- runif(num_p, 0, 1)
```

```
y_p <- runif(num_p, 0, 1)
```

```
plot(x_p, y_p, xlim = c(0, 1), ylim = c(0, 1), xlab = "x", ylab = "y",  
     main = "Plot_of_Parent_Points_for_the_Example")
```

```
loc_par <- ppp(x_p, y_p, window = owin(c(0, 1), c(0, 1)))
```

```
drive_intense <- density(loc_par, sigma = 0.1, kernel = "gaussian")
```

```
plot(drive_intense, main = "Plot_of_the_Driving_Intensity_for_the  
~~~~~Example")
```

```

full_sim <- rpoispp(max(drive_intense), win = owin(c(0, 1), c(0, 1)))

plot(full_sim, main = "Plot_of_the_Full_Homogeneous_Poisson_Simulation")

new_x <- NULL
new_y <- NULL

xvals <- full_sim$x
yvals <- full_sim$y

nvals <- length(xvals)

prob_dens <- as.function(drive_intense / max(drive_intense))

for (i in 1:nvals) {

  prob <- prob_dens(xvals[i], yvals[i])
  ind <- rbinom(1, 1, prob)

  if (ind == 1) {

    new_x <- c(new_x, xvals[i])
    new_y <- c(new_y, yvals[i])

  }
}

thinned <- ppp(new_x, new_y, window = owin(c(0, 1), c(0, 1)))

plot(full_sim, main = "Plot_of_the_Full_Homogeneous_Poisson_Simulation")
points(thinned, pch = 16, col = "red")

```

A.2 Simulation Functions

```

# Description —————

# Script containing the functions to carry out the simulations for the first
# part of the project

# Load Packages —————

library(spatstat)

# Thinning Function —————

# Function to thin a simulation from a homogeneous Poisson process to an
# inhomogeneous Poisson process simulation

# Inputs:
#   intense – intensity given by the density function from spatstat
#   full_sim – ppp object which is a simulation from a homogeneous Poisson
#               with intensity equal to the maximum intensity seen across the
#               window
#   xydim – limit of the window in both the x and y directions (positive number)

# Outputs:
#   ppp object containing the thinned simulation from a Thomas process

thin_sim <- function(intense, full_sim, xydim) {

  # Create variables to store a vector containing the points being kept
  new_x <- NULL
  new_y <- NULL

  # Extract the x and y values into vectors and find the length of the vectors
  xvals <- full_sim$x
  yvals <- full_sim$y

  nvals <- length(xvals)

```



```

# Convert the derived intensity function to the range [0, 1] and then convert
# it to a function type so that the value of the density at (x, y) is
# extracted easily
prob_dens <- as.function(intense / max(intense))

for (i in 1:nvals) {

  # Find probability density at the given point then sample a Bernoulli
  # variable with probability of 1 equal to this probability
  prob <- prob_dens(xvals[i], yvals[i])
  ind <- rbinom(1, 1, prob)

  # If the indicator is 1 then keep the given point and add it to the vector
  if (ind == 1) {

    new_x <- c(new_x, xvals[i])
    new_y <- c(new_y, yvals[i])

  }

}

# Convert to a spatstat ppp object and return
thinned <- ppp(new_x, new_y, window = owin(c(0, xylim), c(0, xylim)))

return(thinned)

}

# Thomas Simulation Function ———

# Function to run the simulation of a Thomas cluster process using the
# thinning function above

# Inputs:

```

```

# mu_p - mean of parent distribution (positive number)
# sd_c - standard deviation of symmetric normal distribution of children
#         around parents (positive number)
# xylin - limit of the window in both the x and y directions (positive number)
#         has a default value of 1)
# rand_seed - number to set the random seed to to allow reproducibility of
#             simulations (number, has an arbitrary default of 150)

# Outputs:
# ppp object containing the simulation from a Thomas process

ThomasSimul <- function(mu_p, sd_c, xylin = 1, rand_seed = 150) {

  # Error traps to ensure that user inputs are valid
  if (!is.numeric(mu_p)) stop("invalid_arguments._mu_p_must_be_a_numeric")
  if (!is.numeric(sd_c)) stop("invalid_arguments._sd_c_must_be_a_numeric")
  if (!is.numeric(xylin)) stop("invalid_arguments._xylin_must_be_a_numeric")
  if (!is.numeric(rand_seed)) stop("invalid_arguments._rand_seed_must_be_a_numeric")

  if (mu_p <= 0) stop("invalid_arguments._mu_p_must_be_>_0")
  if (sd_c <= 0) stop("invalid_arguments._sd_c_must_be_>_0")
  if (xylin <= 0) stop("invalid_arguments._xylin_must_be_>_0")

  # Set random seed to allow reproducibility of simulations
  set.seed(rand_seed)

  # Randomly select number of parents from Poisson distribution of mean mu_p
  num_p <- rpois(1, mu_p)

  # Randomly generate coordinates for each of the num_p points, with equal
  # probability everywhere in the region of [0, xylin] in each direction
  x_p <- runif(num_p, 0, xylin)
  y_p <- runif(num_p, 0, xylin)
  #plot(x_p, y_p, xlim = c(0, xylin), ylim = c(0, xylin))

  # Convert to a spatstat ppp object

```

```

loc_p <- ppp(x_p, y_p, window = owin(c(0, xylim), c(0, xylim)))

# Calculate the driving intensity using Gaussian densities placed at each
# parent location with standard deviation given by sd_c
drive_intense <- density(loc_p, sigma = sd_c, kernel = "gaussian")

# Create sample of a homogeneous Poisson process with density equal to the
# highest seen in the derived driving intensity above
large_sim <- rpoispp(max(drive_intense), win = owin(c(0, xylim), c(0, xylim))

# Thin this using the thinning function described above
thinned_sim <- thin_sim(drive_intense, large_sim, xylim)

return(list(thin_sim = thinned_sim, intense = drive_intense))

}

```