

Spatial Statistics Workshop

Claire Kelling

Pennsylvania State University- Department of Statistics

BIGSSS CSS 2019

Resources: [Shaby, 2017], [Schutte, 2018], [Diggle, 2013], [Schabenberger and Gotway, 2017]

Introduction

Preliminaries

What is spatial data?

- ① Point-referenced (geostatistical)
 - $Y(s); s \in \mathbb{R}^d$ and s varies continuously
- ② Areal (lattice)
 - Finite number of areal units, e.g. counties or elements of a grid
 - Observations are typically sums or averages
- ③ Point Process (point patterns)
 - Locations are themselves the data

Potential Questions of Interest

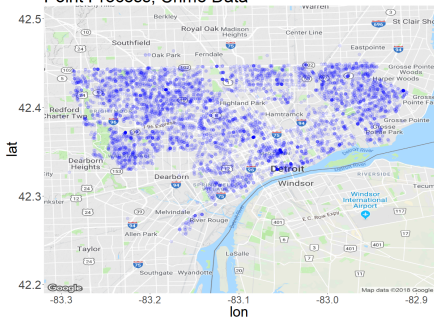
- ① Point-referenced
 - What is the temperature at an un-measured location? (prediction)
- ② Areal
 - Are variations in outcomes related to covariates, such as demographic characteristics? (estimation)
- ③ Point Process
 - Can we distinguish clustering or repulsion?

We will focus on areal unit and point process data.

Spatial Modeling

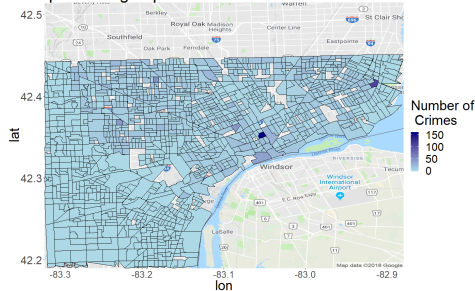
Point Process Modeling

Point Process, Crime Data



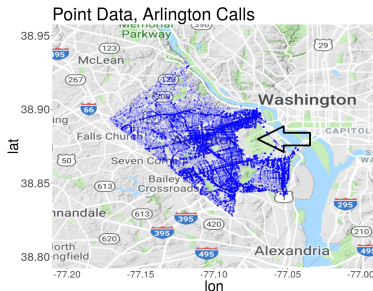
Areal Unit Modeling

Number of Domestic Violence Crimes per block group

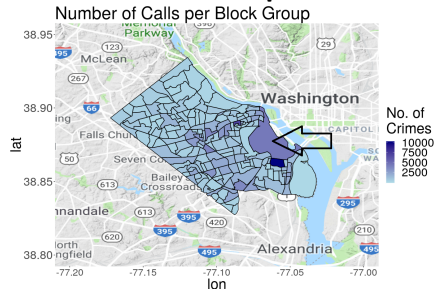


Caution

Crime Point Process



Crime Aggregated by Census Block Groups

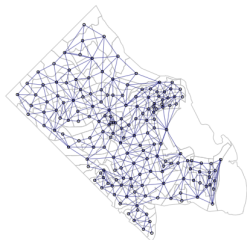


- All of these calls are occurring on the border of Arlington National Cemetery, very close to other communities.
- Also, do not just create maps that are essentially population maps.

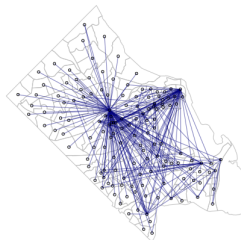
Areal Unit Models

Neighborhood Matrix

Geographic Neighborhood



Social Neighborhood



Goal

Define/capture the dependence between areal units.

Neighborhood Matrix, \mathbf{W}

Some details on neighborhood matrix, \mathbf{W}

- non-negative, symmetric, $K \times K$
- (i,j) th element of the neighborhood matrix w_{ij} represents spatial closeness between areas $(\mathcal{S}_i, \mathcal{S}_j)$
- positive values denoting geographical closeness and zero values denoting non-closeness (0-1 is the most common structure)
- $w_{ii} = 0$

Bayes CAR Model Lee [2013] Set Up

CAR = Conditional Autoregressive model for areal data

- study region \mathcal{S} is partitioned into K non-overlapping areal units
- linked to set of responses $\mathbf{Y} = (Y_1, \dots, Y_K)$
- spatial variation in the response is modeled by a matrix of covariates $\mathbf{X} = (x_1, \dots, x_k)$ and a spatial structure component $\psi = (\psi_1, \dots, \psi_k)$
- $\psi = (\psi_1, \dots, \psi_k)$ models any spatial autocorrelation that remains after covariate effects have been accounted for

GLMM for spatial areal unit data

$$Y_k | \mu_k \sim f(y_k | \mu_k, \nu^2) \text{ for } k = 1, \dots, K$$

$$g(\mu_k) = \mathbf{x}_k^T \beta + \psi_k$$

$$\beta \sim N(\mu_\beta, \Sigma_\beta)$$

- Poisson:

$$Y_k | \mu_k \sim \text{Poisson}(\mu_k) \text{ and } g(\mu_k) = \ln(\mu_k) = \mathbf{x}_k^T \beta + \psi_k$$

BYM Model (Besag, York, and Mollie) [Besag et al., 1991]

$$\psi_k = \phi_k + \theta_k$$

$$\phi_k | \phi_{-k}, \mathbf{W}, \tau^2 \sim N\left(\frac{\sum_{i=1}^K w_{ki} \phi_i}{\sum_{i=1}^K w_{ki}}, \frac{\tau^2}{\sum_{i=1}^K w_{ki}}\right)$$

$$\theta_k \sim N(0, \sigma^2)$$

$$\tau^2, \sigma^2 \sim \text{Inverse-Gamma}(a, b)$$

- First CAR model to be proposed.
- Two sets of random effects, spatially autocorrelated and independent
- Requires two random effects to be estimated at each data point, whereas only their sum is identifiable

Point Process

Point Process Model Set Up

Goal:

To develop a model that will assess for and characterize spatial structure in the data.

- Can develop a flexible model for the mean surface for events, $\mu(s)$, in order to estimate the overall trend across our area of interest.
- Can use a Gaussian Process (GP) model and estimate the parameters of its covariance function.
- Need to assess for Complete Spatial Randomness

Point Process Basics

Poisson Processes over region A

- **Homogeneous Poisson Process:** Static intensity within region A (complete spatial randomness)
- **Inhomogeneous Poisson Process:** Variable intensity within region A
- Depends on *intensity function* $\lambda(s)$

Log Gaussian Cox Process (LGCP)

- $z(s)$ = a vector of location-specific, spatial covariates corresponding to crime event (demographic features, housing information, neighborhood characteristics, etc.)
- $\lambda(s) = r(s)\pi(s)$ where $r(s)$ is the population density at location s , or an offset
- $\pi(s) = \exp(z(s)'\beta + \omega(s))$ where $\omega(s)$ is a zero-centered stochastic process, such as Gaussian Process and β is unknown vector of regression coefficients

We build upon the model developed in Liang et al. [2008] for studying disease risk using spatially varying and non-spatially varying covariates.

Non-spatially varying covariates

$$\pi(s, v) = \exp(\beta_0 + z(s)'\beta + v'\alpha + (v \otimes z(s))'\gamma + \omega(s))$$

where s is the location and v is the crime event

This model allows us to combine covariates at different resolutions, whether it be on crime event or the community level, to draw conclusions about the risk of crime in the area of interest as well as to determine the dominating factors leading to that risk.

Other

Other Topics

Geocoding

- Some of our datasets consist of addresses, not coordinates.
- We have tested various geo-coding techniques to find latitude/longitude coordinates for each of these addresses.
- The results are not always consistent and it can be expensive to get a large number of addresses.

Preliminary Tests

- **Spatial autocorrelation:** Moran's I
- **Complete Spatial Randomness:** Ripley's K

Other Topics

Projections

- It is important to be aware when dealing with spatial data (shape files, point patterns, etc) that projections/coordinate systems are extremely important!
- Projections define what coordinate system you are using (basically)
- Why does this happen?
 - A **map** is a two-dimensional representation of the surface of the earth
 - A **globe** is a three-dimensional model of the earth
 - You lose information when going from 3D to 2D

In R

Packages in R

Popular R packages for spatial statistics

- `maptools`: load and display spatial data
- `rgdal`: reprojections
- `sp`: generate spatial data structures
- `spatstat`: load and manipulate raster data
- `CARBayes`: fit areal unit CAR models
- `ngspatial`: fit modern areal unit models
- `ppm`: fit/explore point process data

Today's Workshop

Find today's workshop presentation and code at

test.com Put on Github or website or both

References/Appendices

References

- Julian Besag, Jeremy York, and Annie Mollié. Bayesian image restoration, with two applications in spatial statistics. *Annals of the Institute of Statistical Mathematics*, 43(1):1–20, 1991.
- Peter J Diggle. *Statistical analysis of spatial and spatio-temporal point patterns*. Chapman and Hall/CRC, 2013.
- Duncan Lee. CARBayes: An R package for Bayesian spatial modeling with conditional autoregressive priors. *Journal of Statistical Software*, 55(13):1–24, 2013.
- Shengde Liang, Bradley P Carlin, and Alan E Gelfand. Analysis of minnesota colon and rectum cancer point patterns with spatial and nonspatial covariate information. *The annals of applied statistics*, 3(3):943, 2008.
- Oliver Schabenberger and Carol A Gotway. *Statistical methods for spatial data analysis*. Chapman and Hall/CRC, 2017.
- Sebastian Schutte. Spatial event data analysis (in r!). *BIGSSS CSS in Conflict*, 2018.
- Ben Shaby. Spatial models. *Introduction to Spatial Statistics*, 2017.

Thanks!

Thanks for your attention!
Questions?

Contact information:

Email: cek32@psu.edu