Multivariate spatiotemporal modeling with applications to stroke mortality and data privacy

Harrison Quick (Drexel University) Joint work with Lance Waller (Emory) and Michele Casper (CDC)

The findings and conclusions in this presentation are mine and do not necessarily represent the official position of the Centers for Disease Control and Prevention.

Table of Contents

Introduction

Methods

Multivariate space-time CAR model Generation and evaluation of synthetic data

Results

Analysis of the stroke mortality data Generation/Evaluation of synthetic data

Summary and Discussion

Table of Contents

Introduction

Methods

Multivariate space-time CAR model Generation and evaluation of synthetic data

Results

Analysis of the stroke mortality data Generation/Evaluation of synthetic data

Summary and Discussion

Goal for this Talk

The charge of agencies such as the CDC includes the following:

- Conduct surveillance into epidemiologic issues
 - e.g., develop/implement statistical models to better estimate and/or predict trends in the data
- Disseminate information (e.g., data) for public use
 - e.g., publishing articles/reports, release data via CDC WONDER
 - Must be cognizant of potential risks of disclosure when sharing information based on confidential/private data

The goal for this talk will be to develop a statistical framework which is useful for both of these charges.

Today's Example: Stroke Mortality

Background information on stroke mortality:

- Stroke is the fourth leading cause of death in the US
- Mortality rates increase exponentially with age
- Previous work has identified strong spatial patterns in stroke mortality (e.g., "the stroke belt")

Our data consists of the number of stroke deaths, Y_{ikt} , and the population size, n_{ikt} , from:

- i = 1,..., N_s=3,099 counties (or county equivalents) from the contiguous United States
- $t = 1, ..., N_t = 41$ years of data (1973 2013)

▶ US citizens ages 65 and older.

▶ $k = 1, ..., N_g = 3$ age brackets (65–74, 75–84, 85+)

Because stroke mortality is quite rare, many of our $N_s \times N_g \times N_t$ = 381,177 counts are quite small.

Data Dissemination Challenges

When releasing these data for public use, CDC WONDER uses NCHS's recommendation of *suppressing* instances where $Y_{ikt} < 10$

 Leads to nearly 70% of the data analyzed here being suppressed.

This has an impact on the types and quality of inference that outside researchers can conduct *using the public-use data*.

- Analyzing all 380,000+ observations would require censored data methods (or otherwise accounting for the missingness)
 — this is likely an unreasonable expectation.
- Others may restrict their analyses to counties in which complete data are available (i.e., urban centers), or aggregate spatially or across age to obtain larger counts.
- Analyses for more specific demographic groups are left unstudied (e.g., mortality rates by age/race/sex), as the issue will only be compounded.

Our Proposal

To obtain more reliable estimates from the data and to provide unrestricted access to high-quality public-use data, we propose the following:

- Analyze the data using a Bayesian statistical model which accounts for (a) spatial structure, (b) temporal structure, and (c) between-age-group structure
 - ► To do so, we will use the multivariate space-time conditional autoregressive (MSTCAR) model of Quick et al. (2017).
- 2. Using the posterior distribution from the Bayesian model, we will generate multiply-imputed *synthetic* data to replace sensitive counts
 - The resulting synthetic data will preserve the complex spatial, temporal, and between-age dependencies (along with any covariate relationships) that we accounted for in our model.

Table of Contents

Introduction

Methods

Multivariate space-time CAR model Generation and evaluation of synthetic data

Results

Analysis of the stroke mortality data Generation/Evaluation of synthetic data

Summary and Discussion

Disease mapping — the univariate case

Following the convention set forth by Besag et al. (1991), we may assume

$$Y_{ikt} \mid \lambda_{ikt} \sim \mathsf{Pois}\left(n_{ikt}\lambda_{ikt}\right)$$
 where $\log \lambda_{ikt} \sim \mathsf{Norm}\left(\mathbf{x}_{ikt}^{\mathsf{T}} \boldsymbol{\beta}_{kt} + Z_{ikt}, \tau_{k}^{2}\right)$,

where

- ★ x^T_{ikt} β denotes a regression where x_{ikt} denotes a vector of county-level covariates
 - For this analysis, our covariates include % non-white and % male within each age group at each time period
- Z_{ikt} denotes a spatiotemporal random effect
- τ_k^2 denotes the variance of the log mortality rates

Conditional autoregressive (CAR) models

To induce spatial correlation in the random effects, Besag et al. (1991) assumed

$$Z_{ikt} | \mathbf{Z}_{(i)kt}, \sigma_{kt}^2 \sim \operatorname{Norm}\left(\sum_{j \sim i} Z_{jkt}/m_i, \sigma_{kt}^2/m_i\right) \\ \pi \left(\mathbf{Z}_{\cdot kt} | \sigma_{kt}^2\right) \propto \left(\sigma_{kt}^2\right)^{-(N_s - 1)/2} \exp\left[-\frac{\mathbf{Z}_{\cdot kt}^{\mathsf{T}} \left(D - W\right) \mathbf{Z}_{\cdot kt}}{2\sigma_{kt}^2}\right]$$

where

- ► $\mathbf{Z}_{(i)kt}$ is the vector $\mathbf{Z}_{kt} = (Z_{1kt}, \dots, Z_{N_skt})^T$ with the *i*th element removed.
- $j \sim i$ denotes that counties i and j are neighbors.
- ► W is an adjacency matrix with w_{ij} = 1 if j ~ i and w_{ij} = 0 otherwise.
 - $m_i = \sum_i w_{ij}$, the number of neighbors
 - D is a diagonal matrix with elements m_i
- σ_{kt}^2 is an age/time-specific variance parameter.

Extension to multiple disease mapping

When modeling data from multiple diseases (or in our case, mortality rates for multiple age groups over time), a multivariate extension of the CAR model can be used (e.g., the multivariate CAR (MCAR) of Gelfand and Vounatsou, 2003).

$$\begin{split} \mathbf{Z}_{i\cdots} \, | \, \mathbf{Z}_{(i)\cdots}, \mathbf{\Sigma}_{Z} &\sim \mathsf{Norm}\left(\sum_{j\sim i} \mathbf{Z}_{j\cdots}/m_{i}, \frac{1}{m_{i}}\mathbf{\Sigma}_{Z}\right) \\ &\pi \left(\mathbf{Z} \, | \, \mathbf{\Sigma}_{Z}\right) \propto |\mathbf{\Sigma}_{Z}|^{-(N_{s}-1)/2} \exp\left[-\frac{1}{2}\mathbf{Z}^{T} \left\{(D-W)\otimes\mathbf{\Sigma}_{Z}^{-1} \mid \mathbf{Z}\right\}, \end{split}$$

where

- ► Z is a N_sN_gN_t × 1 vector of spatiotemporal random effects which allows for correlation between age groups
- Σ_Z is the multivariate analog of σ^2 from the univariate case

Multivariate space-time model for Z

Based on the MCAR of Gelfand and Vounatsou (2003),

$$\mathbf{Z}_{i..} \mid \mathbf{Z}_{(i)..}, \Sigma_Z \sim \mathsf{Norm}\left(\sum_{j \sim i} \mathbf{Z}_{j..}/m_i, \frac{1}{m_i} \Sigma_Z\right)$$

- Spatial associations are accounted for via the neighborhood structure in the mean and variance.
- Thus, Σ_Z can be thought of as a (scaled) covariance matrix which accounts for the multivariate and temporal dependencies in Z.
 - ▶ We'll allow for differing degrees of temporal correlation within each each age-bracket, denoted by $\rho = (\rho_1, \dots, \rho_{N_g})^T$.
 - ▶ Between age-bracket dependencies will be allowed to vary over time, denoted by G = {G₁,..., G_{N_t}}.

We denote this structure by $\mathbf{Z} \sim \mathsf{MSTCAR}(\mathcal{G}, \boldsymbol{\rho})$.

Hierarchical model

Putting these pieces together, our full hierarchical model is as follows:

$$\begin{aligned} \pi \left(\boldsymbol{\beta}, \boldsymbol{\mathsf{Z}}, \boldsymbol{\mathsf{G}}, \boldsymbol{\mathcal{G}}, \boldsymbol{\rho}, \left\{ \tau_k^2 \right\}, \boldsymbol{\lambda} \, | \, \boldsymbol{\mathsf{Y}} \right) & \propto \prod_{i,k,t} \mathsf{Pois} \left(\boldsymbol{Y}_{ikt} \, | \, \boldsymbol{n}_{ikt} \boldsymbol{\lambda}_{ikt} \right) \\ & \times \prod_{i,k,t} \mathsf{Norm} \left(\log \lambda_{ikt} \, | \, \boldsymbol{\mathsf{x}}_{ikt}^T \boldsymbol{\beta}_{ikt} + \boldsymbol{Z}_{ikt}, \tau_k^2 \right) \\ & \times \mathsf{MSTCAR} \left(\boldsymbol{\mathsf{Z}} \, | \, \boldsymbol{\mathcal{G}}, \boldsymbol{\rho} \right) \times \mathsf{Norm} \left(\boldsymbol{\beta} \, | \, \boldsymbol{\mathsf{0}}, \boldsymbol{\boldsymbol{\Sigma}}_{\boldsymbol{\beta}} \right) \\ & \times \prod_t \mathsf{InvWish} \left(\boldsymbol{\mathsf{G}}_t \, | \, \boldsymbol{\mathsf{G}}, \boldsymbol{\nu} \right) \times \mathsf{Wish} \left(\boldsymbol{\mathsf{G}} \, | \, \boldsymbol{\mathsf{G}}_0, \boldsymbol{\nu}_0 \right) \\ & \times \prod_k \left[\mathsf{Beta} \left(\rho_k \, | \, \boldsymbol{a}_{\rho}, \boldsymbol{b}_{\rho} \right) \times \mathsf{IG} \left(\tau_k^2 \, | \, \boldsymbol{a}_{\tau}, \boldsymbol{b}_{\tau} \right) \right], \end{aligned}$$

where $\Sigma_{\beta} = 100 I_{pN_gN_t}$ and X is the $(N_s N_g N_t \times p)$ matrix of covariates.

We fit this model using Markov chain Monte Carlo (MCMC) and obtain samples from the posterior distribution for each model parameter.

• e.g., $\lambda_{ikt}^{(1)}, \ldots, \lambda_{ikt}^{(L)}$, where L is the number of iterations

Synthetic data

Given our samples for λ_{ikt} , we can generate synthetic counts for our suppressed Y_{ikt} from a truncated Poisson of the form

$$Y_{ikt}^{*(\ell)} \,|\, \lambda_{ikt}^{(\ell)}, \{ Y_{ikt} < 10 \} \sim \mathsf{Pois}\left(\mathsf{n}_{ikt} \lambda_{ikt}^{(\ell)} \right) \times \mathsf{I}\left\{ Y_{ikt}^{*(\ell)} < 10 \right\}.$$

If desired, this approach could be modified to preserve aggregate totals (e.g., state-level counts) which would be publicly available.

To assess the quality of these synthetic data, we will compare them to synthetic data that could be generated by fitting the MSTCAR model to the publicly available (i.e., *suppressed*) data.

- Counts below 10 will be imputed as part of the model
- We consider this to be the best available alternative for both public users and for ill-intentioned users (or "intruders")

Measuring disclosure risk and utility

Disclosure risk will be computed as

$$P(Y_{ikt}^* = y | \mathbf{Y}, Y_{ikt} = y)$$
 for $y = 0, 1, \dots, 9$.

In particular, we will look at the risk when y = 1 (the value we're most concerned about).

Utility will be compared by fitting a model of the form

$$Y_{ikt} \sim \mathsf{Pois}\left(n_{ikt}\exp\left[\gamma_{0kt} + \mathsf{rural}_{ikt}\gamma_{1kt}
ight]
ight),$$

where rural_{*ikt*} denotes a 0/1 variable taking value 1 if county *i* has a population (across all age groups) less than 50,000 during year *t*.

Estimates from synthetic data will also be compared to the estimates from the confidential data (i.e., the "truth").

Table of Contents

Introduction

Methods

Multivariate space-time CAR model Generation and evaluation of synthetic data

Results

Analysis of the stroke mortality data Generation/Evaluation of synthetic data

Summary and Discussion

Stroke mortality: ages 65-74

1973



Overall declines in stroke mortality



(c) Ages 85+

Over 7756

How much of these data are suppressed to the public?



Example: 1986* in Montour County, PA



* Data *since* 1989 is suppressed on CDC Wonder, but data *prior* to 1989 is unsuppressed and publicly available.

Disclosure risk



(a) $P(Y_{ikt}^* = 0 | Y_{ikt} = 0)$ (b) $P(Y_{ikt}^* = 1 | Y_{ikt} = 1)$ (c) $P(Y_{ikt}^* = 9 | Y_{ikt} = 9)$

- Red and green lines denote the expected risk probabilities at the beginning and end of the study, respectively.
- These risk probabilities are highest at the boundary values.
 - If $Y_{ikt} = 0$, there is no one's privacy to be concerned about.
 - We set the upper bound to some conservative value.
- Interior values are essentially what we would "expect"

Disclosure Risk and Utility



Table of Contents

Introduction

Methods

Multivariate space-time CAR model Generation and evaluation of synthetic data

Results

Analysis of the stroke mortality data Generation/Evaluation of synthetic data

Summary and Discussion

Summary

Recall that the goal of this talk was to develop a statistical framework which is useful for both public health surveillance and the dissemination of information, thereby avoiding a redundancy of tasks. Thus, we claim:

- The MSTCAR is well-suited for conducting public health surveillance.
 - The posterior distribution yields inference on rates, aggregates of rates, rate ratios, declines, etc.
- The MSTCAR shows promise for generating synthetic data for public-use
 - Using the MSTCAR should yield synthetic data with very high utility
 - That said, it is not without its weaknesses

Limitations / Future Work

- No clear connection (yet) between this approach and a form of differential privacy
 - ► We see some similarities between our framework and that used for OnTheMap (Machanavajjhala et al., 2008), but the question is how to express the "informativeness" of our model.
- Not practical for BIG examples without BIG assumptions
 - A similar analysis with $N_g = 24$ age/race/sex subgroups takes 2+ weeks to run
- Aspects of utility unclear
 - e.g., we assume (but haven't proven) that by accounting for spatial structure, we will preserve relationships for spatially-structured covariates *not included* in the model

Our vision: For this approach to ultimately be used for a series of one-offs rather than to generate a "Synthetic CDC WONDER"

 e.g., CDC researchers study trends in stroke mortality, publish their research, and make the synthetic data available for further analysis by outside researchers

Questions?

hsq23@drexel.edu