



# Bayesian hierarchical modeling of the dynamics of spatio-temporal influenza season outbreaks

Andrew B. Lawson<sup>a,\*</sup>, Hae-Ryoung Song<sup>b</sup>

<sup>a</sup> Division of Biostatistics & Epidemiology, Dept of Medicine, Medical University of South Carolina, Charleston, USA

<sup>b</sup> Department of Research Affairs, Yonsei University, Seoul, Republic of Korea

## ARTICLE INFO

### Keywords:

Influenza

Modeling

SIR

C positive

Under-ascertainment

Bayesian

## ABSTRACT

The analysis of influenza incidence in space and time is considered. A SIR model is proposed for the bi-weekly C+ notifications of lab confirmed influenza for the 2005 flu season. A variety of models are considered and the resulting goodness of fit and other diagnostics are considered.

© 2010 Elsevier Inc. All rights reserved.

## 1. Introduction

The modeling of infectious disease has seen a considerable development both in terms of etiological understanding (Anderson and May, 1991; Fraser et al., 2004) as well as theoretical development (see e.g. Daley and Gani, 2000; Andersson and Britton, 2000; Ball and Britton, 2009). Recently individual level space–time models have been developed in the particular case of measles outbreaks (Lawson and Leimich, 2000; Neal and Roberts, 2004). For aggregate count data, a variety of temporal models have been proposed (see e.g. Finkenstädt et al., 2002; Morton and Finkenstädt, 2005; Cauchemez and Ferguson, 2008). Morton and Finkenstädt (2005) have modeled the progression of measles outbreaks over time using a susceptible–Infected–removed (SIR) model. For influenza, there are few examples of attempts to use standard public health data as a basis for modeling (Finkenstädt et al., 2005). In this paper we investigate the possibility of developing space–time models for publicly available influenza data at the county level in the US state of South Carolina. We extend the basic model suggested by Lawson (2006) by including correlated and uncorrelated random effects, and also suggest a new model of infection process by con-

sidering neighborhood infection effects. The paper is divided as follows. First we introduce the available data and discuss its limitations, in particular under-ascertainment. Next, we consider possible SIR-type models and their extensions into spatial data. We then discuss computational issues and model fitting. Finally, we consider future developments.

## 2. Data example

Our focus in this study is the analysis of publicly available influenza data from the acute influenza data archive of the South Carolina Department of Health and Environmental Control (SCDHEC). The online system can be searched at <http://www.scdhec.gov/health/disease/acute/flu.htm> where a variety of data sources are available. In our example and modeling, we focus on the culture positive (C+) laboratory notifications of influenza which are available by county for the 46 counties of South Carolina (SC). The data is available as counts of C+ notifications for bi-weekly periods at the county level. These notifications represent cases of confirmed influenza within a population whose total influenza burden is much greater. The under-ascertainment for this data is considerable. However, the general pattern of confirmations is likely to mirror the overall behavior of the epidemic progression and so we can consider these notifications as a form of sentinel or syndromic variable

\* Corresponding author.

E-mail address: [lawsonab@musc.edu](mailto:lawsonab@musc.edu) (A.B. Lawson).

useful in the prediction of the epidemic events. We examine the main influenza season from October 2004 through to April 2005 with 13 bi-weekly time periods. The data is available for successive flu seasons, but for our purposes we examine a single season in this paper (October 18th 2004 to April 30th 2005). Fig. 1 displays the selection of crude count time profiles for 4 counties representing some of the main population centers in the state. Fig. 2 displays a selection of cumulative crude count maps for 4 time periods: week beginning 15th December 2004, 15th January, 1st March and 1st April 2005. Flu season definition is an issue and the time of start and finish could be considered to be an important indicator. However, we do not consider this estimation problem in his study. Our modeling assumes a fixed time period during which influenza outbreaks occur.

We also only examine a single season. However, our methods can be extended to multiple seasons without difficulty.

### 3. Model development

As the data are derived as counts of new cases, we will consider these as the main data source in this study. Denote the count of C+ notifications in the  $i$ th county and at the  $j$ th time period as  $y_{ij}$ ,  $i = 1, \dots, m$ ;  $j = 1, \dots, J$  where  $m = 46$  and  $J = 13$ . Our focus is on modeling infective counts ( $I_{ij}$ ) in the counties of SC and so we must provide some model link between laboratory notifications (C+) and total

infectives. Infectives are assumed to be infected individuals.

#### 3.1. C+ notifications (CPNs) and infectives

The relation between observed CPNs and infectives can be specified in a variety of ways. The simplest choice for a model is the proportional link where the total infectives are assumed to be proportional to CPN. A fixed proportion can be assumed thus:

$$y_{ij} = \rho I_{ij} \quad (1)$$

The fixed rate assumed implies a homogeneous non-stratified relation between CPNs and counts of new infectives in a given region. Note that this idea can be extended to allow variation between regions by assuming a distribution for this relationship. A simple case would be a binomial distributional assumption where:

$$y_{ij} \sim \text{bin}(\rho, I_{ij}). \quad (2)$$

This is a very simple way of dealing with the estimation of 'true' infectives. If stratification of the infective counts was available within each county then more complex stratified relations could be assumed. We do not pursue this here.

#### 3.2. The infective model

Our model for incident cases of disease within each county assumes first that the probability distribution gov-

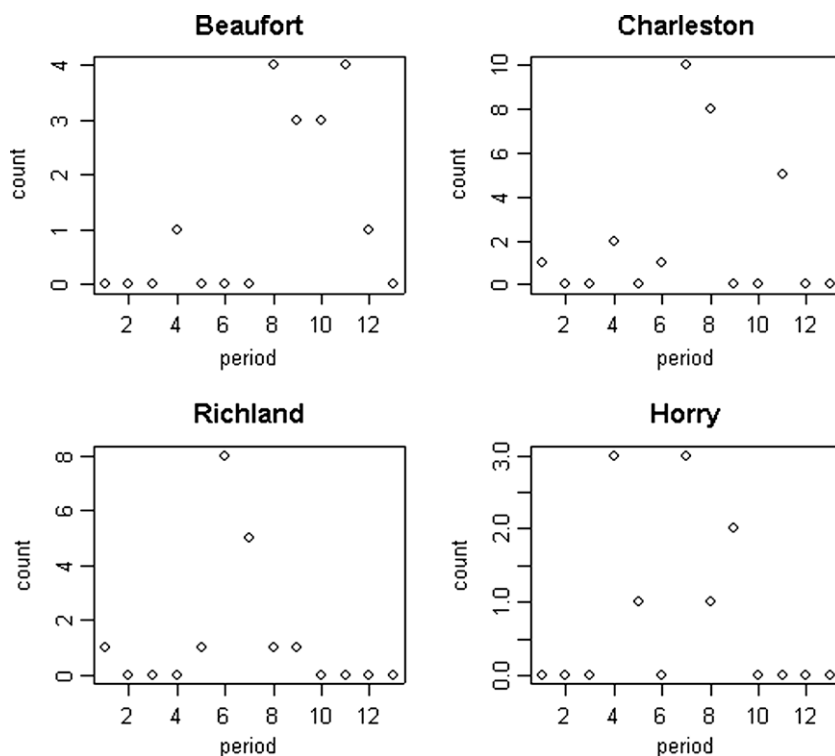
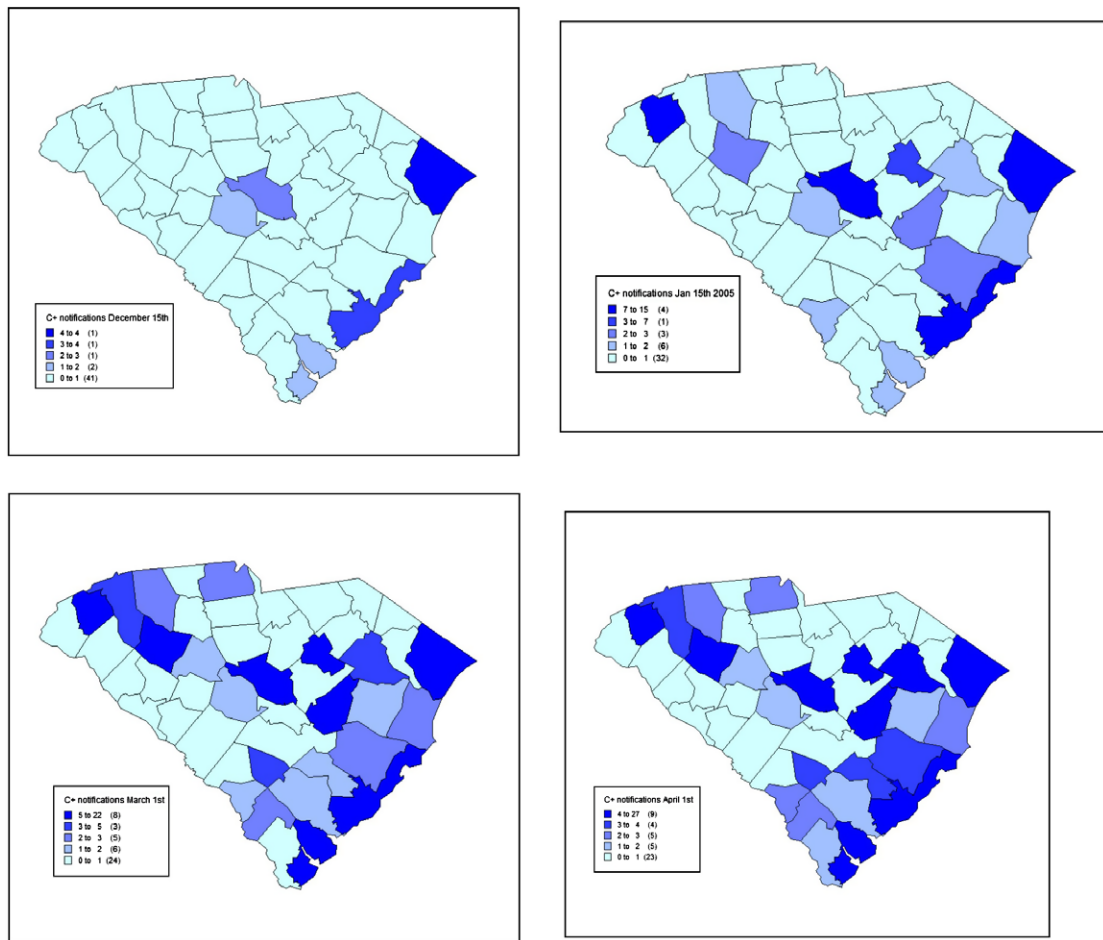


Fig. 1. Time profiles of bi-weekly counts for four counties of South Carolina: Beaufort, Charleston, Richland and Horry during the 2004/2005 influenza season.



**Fig. 2.** A selection of three county maps for fixed time periods during the flu season of 2004/2005: cumulative counts week beginning 15th December 2004, 15th January, 1st March and 1st April 2005.

erning the counts is Poisson and that we can assume conditional independence of counts given complete knowledge of the hierarchical effects within the influenza syndrome. We assume that

$$I_{ij} \sim \text{Pois}(\mu_{ij}) \quad (3)$$

Hence, a Poisson likelihood can be assumed where

$$L(\mathbf{I}_{ij}|\theta) = \prod_{i=1}^I [\mu_{ij}^{I_{ij}} \exp\{-\mu_{ij}\} / I_{ij}!]$$

where  $\mu_{ij} = f(\theta; \mathbf{I}_{ij-1})$ .

Temporal dependence is modeled within the mean level of the process and a variety of possible models can be specified following this assumption.

### 3.3. Observable influenza dynamics

With only the incident count of CPNs available bi-weekly in each county, we do not have information about movement of infected people nor detail of contact distributions, and only partial information about total infective burden in each county. If we knew contact distri-

butions or had migration information at county level over bi-weekly periods we could build such transmission routes into our formulation. Instead we must resort to the use of temporal and spatial dependency to make allowance for these effects.

Within this largely rural state there is little long range movement except possibly between the main urban centers and on interstate highways transiting the state.

At the county level the simplest approach to this largely local effect is to assume that neighboring counties could affect transmission of disease. Note that the definition of neighbor is crucial here as we could define neighborhoods that are adjacent to each other or even not adjacent so as to allow jump diffusions between distant locations. An alternative way to describe spatial transmission could be to consider distance-based effects where a continuous measure of association is based in distance decline so that, for example,  $(1/\text{distance})$  or  $(\exp(-\text{distance}))$  are used as measures of association. Given the crude level of spatial aggregation in our data we assume that a simpler approach based on an adjacency neighborhood should adequately describe the association.

### 3.4. The SIR model

A susceptible–infected–removed (SIR) model is assumed for the influenza dynamics. This model is an extension of Morton and Finkenstädt (2005) into the spatial domain. Their application was bi-weekly measles notifications, but the basic modeling strategy remains valid in general for influenza.

We consider the following model components. First, the population at risk of influenza forms a susceptible population ( $S_{ij}$ ) and second a proportion of the population who have been infected are removed from the overall population of the study area. Re-infections are rare events during a single influenza season. Therefore, the SIR model considers the “R” component as “removed” and not “recovered”, i.e. recovered individuals are considered as immune and thus do not re-enter the susceptible component. Finkenstädt et al. (2005) consider additional antigenic drift. We do not pursue that here. Note also that the susceptible pool is relatively large in the case of influenza within a general population and so a finite population model is not considered necessary in this case.

### 3.5. Accounting equation

The basic accounting equation states that, at a given time, a new susceptible population is produced by the following balance:

$$S_{ij+1} = S_{ij} - I_{ij} - R_{ij}. \quad (4)$$

$R_{ij} = \beta I_{ij}$  where  $R_{ij}$  is the count of removed cases in the  $i$ th county and at the  $j$ th time period and  $\beta$  is the removal rate from infected cases.

We extend the basic definition of the mean function for infectives to include dependence on the reservoir of susceptibles:

$$\mu_{ij} = f(\theta; \mathbf{I}_{ij-1}; S_{ij})$$

where  $\theta$  is a parameter vector which can include transmission rate and other parameters.

### 3.6. Model variants

A basic mean level model can be proposed where

$$\mu_{ij} = S_{ij} \cdot g(I_{ij-1})$$

where  $g(I_{ij-1})$  takes various forms. First, on the log scale we can assume a linear form:  $\log(\mu_{ij}) = \log(S_{ij}) + \log(g(I_{ij-1}))$  and hence parameterize  $\log(g(I_{ij-1}))$ .

For example, a simple model could be

$$\log g(I_{ij-1}) = \log I_{ij-1} + b_0 + b_i \quad (5)$$

where  $\exp(b_0)$  is the transmission rate and  $\exp(b_i)$  can be considered to be a random effect at the county level. This latter effect could be useful in absorbing extra variation due to unobserved confounders in the data. This simple approach could be further extended by assuming a more complex random effect term:  $b_i = v_i + u_i$  where  $v_i + u_i$  combines to absorb the uncorrelated and spatially correlated random effects due to confounding. This latter com-

ponent is sometimes known as a convolution model (see e.g. Lawson (2009), chapter 5).

Alternative specifications which allow more extensive spatial infective dependencies can be imagined. The model in (5) only depends on infectives in the same area at the previous time. We could extend the dependency in time to further lags e.g.

$$\log g(I_{ij-1}) = \log[I_{ij-1} + I_{ij-2} + \dots] + \text{other terms}. \quad (6)$$

Additionally, we could imagine that neighboring areas could contribute infectives to the current pool (either lagged in time or not). A variety of forms could be considered. Here we will consider a simple model where neighbors at a previous time inform the current state:

$$\log g(I_{ij-1}) = \log \left[ I_{ij-1} + \sum_{k=1}^{n_{\delta_i}} I_{k,j-1} \right] + \text{other terms}. \quad (7)$$

where  $n_{\delta_i}$  is the number of neighboring counties around the  $i$ th county. More complex variants could be easily imagined where more extensive dependencies arise e.g.

$$\log g(I_{ij-1}) = \log \left[ \sum_{k=1}^{n_{\delta_i}} w_k I_{k,j-1} + \sum_{k=1}^{n_{\delta_i}} w_k I_{k,j-2} \right] + \text{other terms}. \quad (8)$$

where  $\{w_k\}$  is a set of neighborhood weights allowing different dependencies on different neighbors. Further these could be varied over time allowing changes in dependencies:  $\{w_{k,j}\}$ .

In each of the models (6)–(8) above the ‘other terms’ is meant to represent additional effects that we want to include to control for either confounding or over dispersion or for the incorporation of covariates. For example, we could have a transmission rate parameter as an intercept plus random effects:  $b_0 + v_i + u_i$ , or alternatively a covariate model  $b_0 + x_i^T \mathbf{b}$  where  $\mathbf{b}$  is a parameter vector, or a mix of these terms. In our data example we do not have covariates readily available and so we confine our additional effects to random effect terms.

## 4. Application

### 4.1. Bayesian model posterior sampling

Our models developed are Bayesian and hence we have examined posterior sampling algorithms via MCMC for parameter estimation. The joint models specified with Poisson data likelihoods and accounting equations are relatively straightforward to specify as hierarchical models. For simplicity, the removal rate (beta) is assumed to be fixed to 0.001 and we assume prior distributions for all relevant stochastic parameters, including  $b_0$ ,  $v_i$ ,  $u_i$  and hyper-parameters in their prior distributions. Specifically, we assume a non-informative Gaussian prior distribution with large variance for  $b_0$ , and uniform hyper-prior distributions for standard deviations related to the precisions of the zero-mean Gaussian and CAR distributions for random effects:  $\tau_* = 1/\sigma_*^2$ ;  $\sigma_* \sim U(0, c)$  (Gelman, 2006). All other parameters are assumed fixed in the analysis reported here. In our reported evaluation of models we have used the

deviance information criterion (DIC) and effective number of parameters (pD) as measures (Spiegelhalter et al., 2002; Lawson, 2009, chapter 4). The primary measure used was DIC which is adjusted for the number of parameters in the model. A secondary measure which can be used to judge parsimony for models that have competitive DICs is the effective number of parameters. For both DIC and pD a lower value is judged to represent an improved model.

For posterior sampling from the joint models we have used the WinBUGS software package (Lunn et al., 2009). Our evaluation of models has considered both direct modelling of CPNs with fixed estimates of total infectives, and also direct modeling of unobserved latent infectives. We have found that direct modeling of infectives as latent variables in space–time is relatively unstable and results were not convergent as opposed to CPNs. Instead we have found that direct modeling of CPNs has been successful in terms of convergence of sampling and these models also provide total infective estimates as posterior functionals of the mean CPNs. Hence we have applied models 4, 5, and 7, in particular, with CPN replacing infective, and thence estimation of total infective from the relation  $I_{ij}^g = \rho^{-1} \mu_{ij}^g$  where  $g$  denotes the  $g$ th iteration and  $\mu_{ij}^g$  is the CPN mean parameter at the  $g$ th iteration. For South Carolina, the average  $\rho$  parameter is known to be 0.0003525 (<http://www.scdhec.gov/health/disease/acute/flu.htm>). Hence the posterior average of  $I_{ij}^g$  will provide a reasonable average estimate over all the counties. Table 1 displays the results of fitting these models to the bi-weekly flu season for 2004–2005. Appendix 1 displays the WinBUGS code for the model with Dep1 with a CAR and UH component in the transmission model (3rd model in Table 1).

It is apparent from Table 1, for the models fitted, and based on the DIC criterion, that individual county lagged time dependence is generally a better fit than a lagged neighborhood model with any combination of random effects. In either category, the lowest DIC model is the model with only an uncorrelated random effect. In general, the convolution or single CAR models do not yield as good models. The UH models are also parsimonious in that they have the lowest pD among random effect models. The best fitting model overall yields a transmission rate posterior mean estimate of  $\exp(-11.52) = 9.929504e-06$ , with random effect precision 1.286 (sd: 0.558). For the best model we have selected a range of posterior average output. Fig. 3 displays the posterior mean average CPN temporal profiles with associated credible intervals for 4 major

urban regions of the state: Beaufort, Charleston, Horry, and Richland. These profiles demonstrate that considerable asynchronous epidemic waves occur during a given flu season, with Horry county signaling first at 4–5th period (mid November–mid December), along with a small peak in Beaufort. Following this the main peaks occur first in Richland at the 6–8th period (early January–late January) and Charleston at the 8–10th period (late January–late February). Horry county also signals during this later period followed by Beaufort which has its major wave between the 9th and 12th period (early February–late March). Charleston also has subsidiary wave in late March.

Figs. 4 and 5 display, for a selection of time periods the posterior average mean CPN and infective maps for the counties of SC. It is noticeable that considerable spatial differentiation occurs during the season with different areas signaling at different times. Note that the infective map may under-estimate the true infective numbers when zero CPNs are observed. Fig. 6 displays the overall uncorrelated random effect at county level. It is notable that for the best model in terms of goodness of fit (DIC), the uncorrelated effect displays some evidence of residual structure. In particular, the area in the south west including the highest effects in Bamberg and Allendale counties suggests some clustering of risk unexplained by this model. Fig. 7 displays the posterior average mean CPN profiles for the residual elevated counties (Allendale, Bamberg and Lee). It is marked that they appear to have short lived but marked peaks of infection mainly in period 8–9 in the west (Allendale and Bamberg) and around period 10–12 in the east (Lee). These counties are largely rural and some of the observed behavior may be ascribable to the isolation effects. It may be supposed that the form of these peaks is not well modeled by the current assumed model, and this is an issue which will demand further investigation.

## 5. Discussion and conclusions

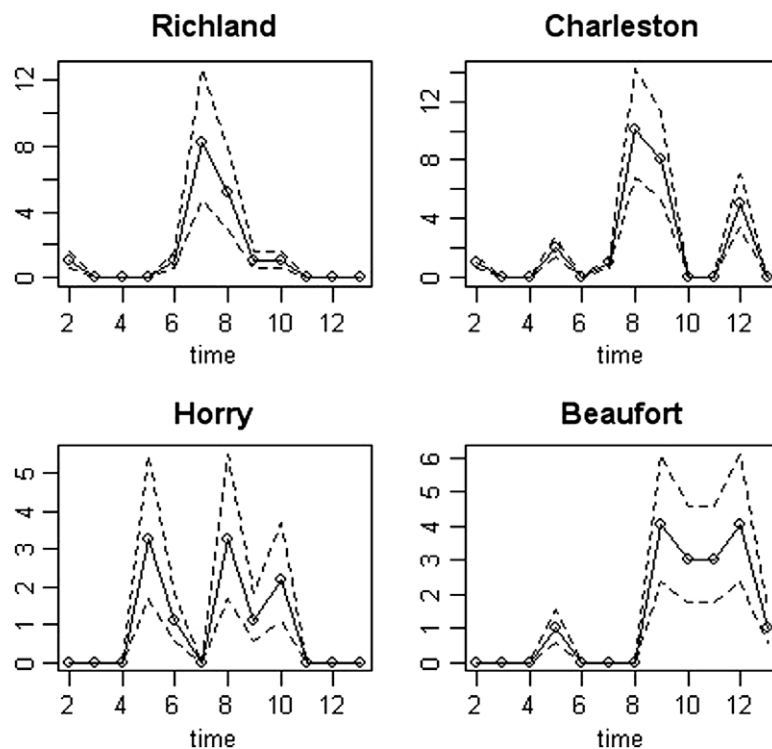
In this paper we have examined the use of Bayesian space–time models for the estimation of infection risk during the influenza season in South Carolina. Our approach has discovered distinct space–time differentiation of the epidemic peaks across the state. There certainly is differential timing of incipient epidemic waves and also differences in duration. Some of these differences could be explained by population characteristics. For example, we might expect areas with high retiree population numbers

**Table 1**

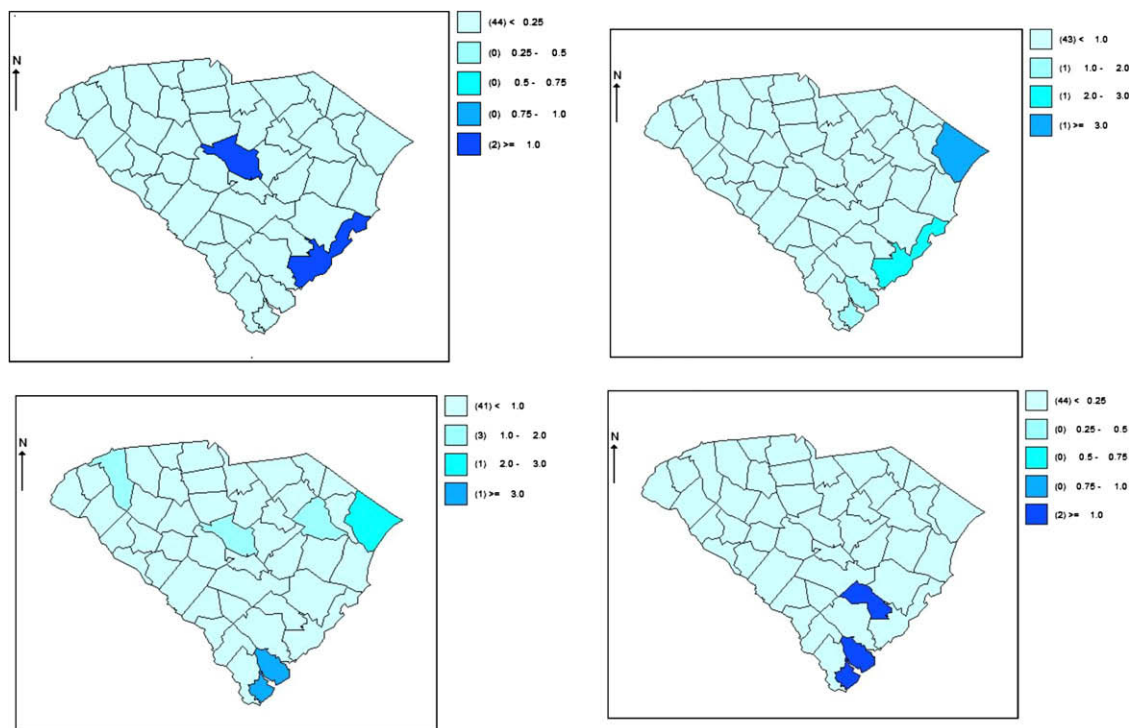
Bayesian SIR model results for CPN as modeled outcome: goodness of fit based on DIC and pD and parameter estimates. Upper panel are models with a single temporal lag dependence on previous time in same county (Dep1). Lower panel are models where a lagged 1 neighborhood of adjacent counties is assumed (Dep neighborhood).

Model	DIC	pD	$b_0$ (sd)	Random effect	Precision (sd)
Cpos: Dep1	8025310	19.9	–11.55 (0.252)	CAR only	0.3585 (0.1504)
	8025310	15.52	–11.52 (0.237)	UH only	1.286 (0.558)
	8025310	17.06	–11.46 (0.320)	UH + CAR	CAR: 5.014 (10.66); UH: 39.43 (217.3)
	8025390	1.03	–12.05 (0.086)	No RE	NA
Cpos: Dep neighborhood	8026130	45.05	–11.06 (0.512)	CAR only	0.0648 (0.0228)
	8026090	21.60	–11.22 (0.544)	UH only	0.1983 (0.0689)
	8026090	22.5	–12.05 (1.655)	UH + CAR	CAR: 1.287 (3.627) UH: 0.192 (0.152)
	8026400	0.98	–12.42 (0.085)	No RE	NA





**Fig. 3.** Posterior mean average CPN infection levels for four urban counties of South Carolina: Richland (Columbia), Charleston, Horry (Myrtle Beach) and Beaufort. Mean levels (solid lines) with associated 95% credible intervals (dotted lines).



**Fig. 4.** Posterior mean average CPN for the UH only lag one model for 4 time periods: row-wise from top left: period 2 (early November), 5 (mid December), 10 (late February), 13 (late April).

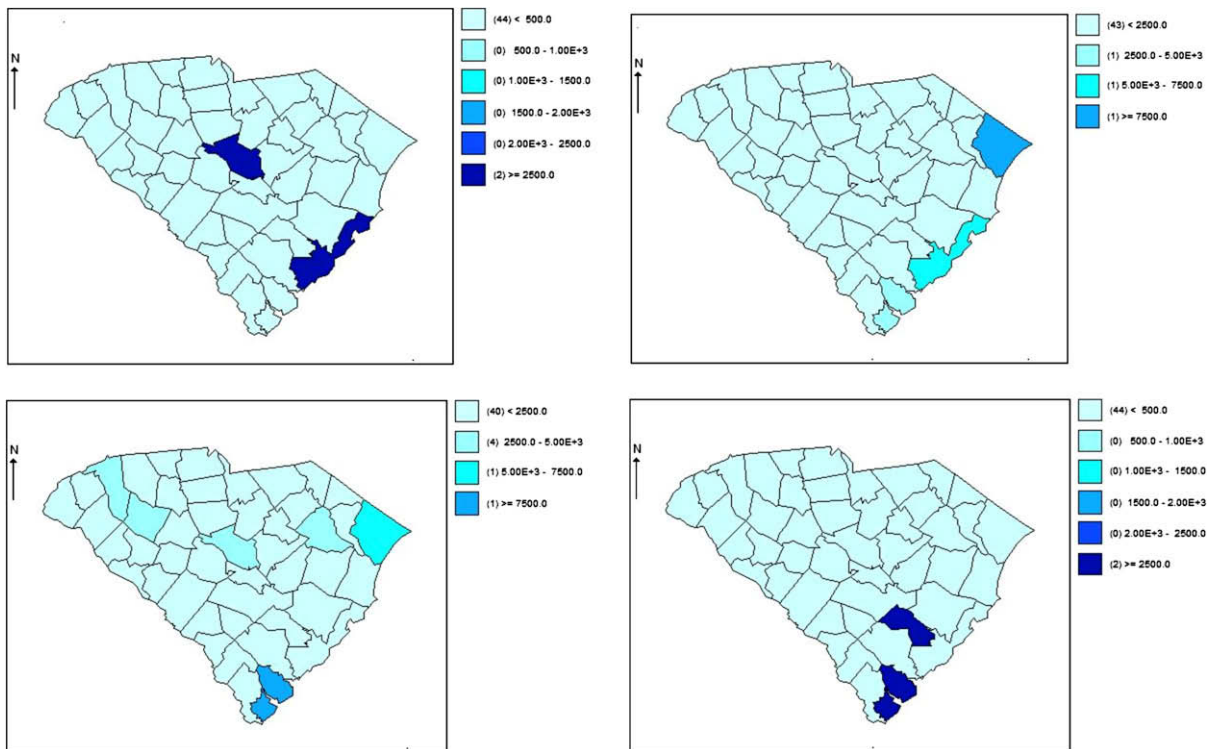


Fig. 5. Posterior mean average Infective numbers for 4 time periods during the 2004–2005 influenza season (periods as for Fig. 4).

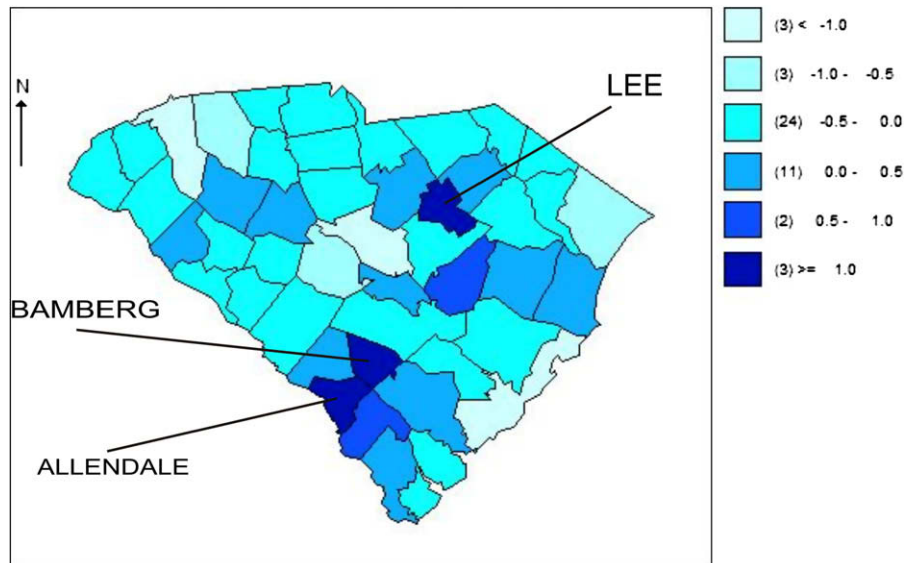


Fig. 6. The spatial distribution of the county-level posterior mean uncorrelated heterogeneity effect: elevated counties identified.

to be affected early and it is true that Horry county, including the beach resort of Myrtle Beach which is popular with retirees, is affected across the time periods whereas Beaufort (including Hilton Head Island) is affected only in the late stages of the season.

In this analysis we have made a variety of assumptions which may be important in determining estimates

of risk and influenza spread. Of these the assumptions, that of uniform under-ascertainment of infectives over time and space may be questioned. Hence Fig. 5 may be a simplification of the link of notification to larger outbreak. In addition we have assumed a constant removal rate and this may be an assumption that could be challenged. Both differential under-ascertainment,

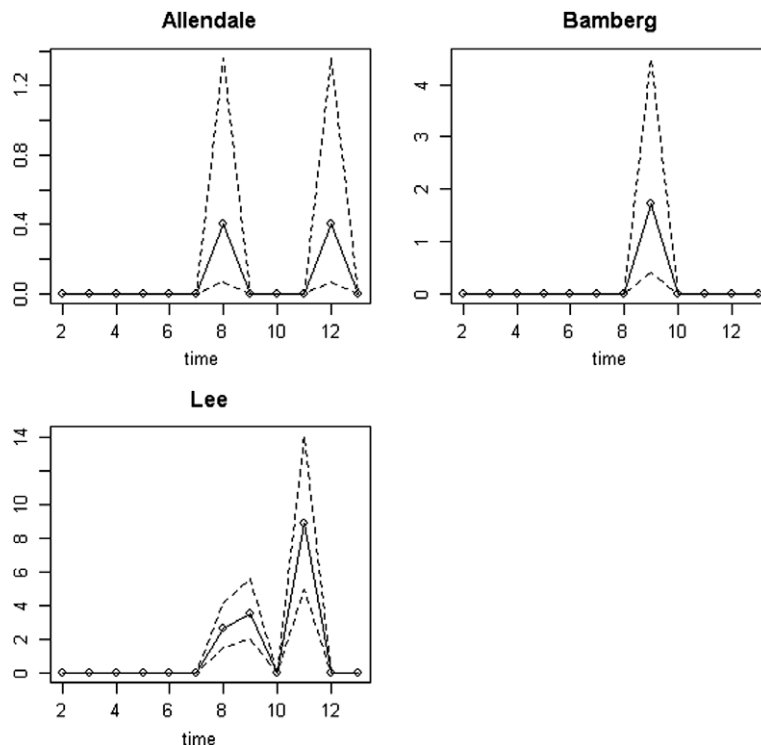


Fig. 7. Posterior average mean CPN temporal profiles for the 3 counties: Allendale, Bamberg, and Lee.

various different lag models and also removal variation (and hence antigenic drift) will be addressed in future work.

Another issue is sensitivity to prior specification and in particular assumptions about random effects in models. In our experience, in examining ranges of different models with different assumptions for random effect hyper-priors we have found that choice of a uniform hyper-prior on the standard deviation of Gaussian random effects is a robust specification.

We do believe that the approach advocated here has practical application in the modeling and prediction of epidemic behavior for routinely collected influenza notification data, and we would hope that given the accessibility of the software a wider use for these models will be found in the public health community.

## Appendix A

As an example of the model code that can be developed for these data we here display the WinBUGS code for the model with CAR and UH component in the transmission model.

```
model{
  for (i in 1:M){
    cumN1[i]<-cumN[i]
    rem[i,1]<-0
    susc[i,1]<-susint[i]
    muc[i,1]<-susc[i,1]
```

```
    cpos[i,1]~dpois(muc[i,1])
    lnf[i,1]<-muc[i,1]/0.0003525
  }
  for (i in 1:M){
    for (j in 2:T){
      rem[i,j]<-betaR*cpos[i,j]
      susc[i,j]<-susc[i,j-1]-cpos[i,j-1]-rem[i,j-1]
      cpos[i,j]~dpois(muc[i,j])
      lnf[i,j]<-muc[i,j]/0.0003525
      log(muc[i,j])<-bet0+log(susc[i,j]+0.001)+
        log(cpos[i,j-1]+0.001)+b1[i]+b2[i]
    }
    b2[i]~dnorm(0,tau.b2)
  }
  for (j in 1:T){
    mucrich[j]<-muc[40,j]
    mucchar[j]<-muc[10,j]
    muchor[j]<-muc[26,j]
    mucbea[j]<-muc[7,j]
    b1[1:46] ~ car.normal(adj[], weights[], num[], tau.b1)
    for(k in 1:sumNumNeigh){
      weights[k] <- 1
    }
    bet0~dflat()
    tau.b2<-pow(sdb2,-2)
    sdb2~dunif(0,5)
    tau.b1<-pow(sdb1,-2)
    sdb1~dunif(0,5)
    betaR<-0.001
  }
}
```



## References

- Anderson RM, May RH. Infectious diseases of humans: dynamics and control. Oxford and New York: Oxford University Press; 1991.
- Andersson H, Britton T. Stochastic epidemic models and their statistical analysis. Springer Lecture Notes in Statistics 2000;vol. 151. New York: Springer-Verlag; 2000.
- Ball FG, Britton T. An epidemic model with infector and exposure dependent severity. *Math Biosci* 2009;218:105–20.
- Cauchemez S, Ferguson N. Likelihood-based estimation of continuous-time epidemic models from time-series data: application to measles transmission in London. *J R Soc Interface* 2008;5:885–97.
- Daley D, Gani J. Epidemic modelling: an introduction. New York: Cambridge University Press; 2000.
- Finkenstädt BF, Bjørnstad ON, Grenfell BT. A stochastic model for extinction and recurrence of epidemics: estimation and inference for measles outbreaks. *Biostatistics* 2002;2002(3): 493–510.
- Finkenstädt BF, Morton AR, Rand DA. Modelling antigenic drift in weekly flu incidence. *Stat Med* 2005;24:3447–61.
- Fraser C, Riley S, Anderson RM, Ferguson NM. Factors that make an infectious disease outbreak controllable. *Proc Natl Acad Sci USA* 2004;101(16):6146–51.
- Gelman A. Prior distributions for variance parameters in Hierarchical models. *Bayesian Anal* 2006;1:515–53.
- Lawson AB. Statistical methods in spatial epidemiology. 2nd ed. New York: Wiley; 2006.
- Lawson AB. Bayesian disease mapping: hierarchical modeling in spatial epidemiology. New York: CRC Press; 2009.
- Lawson AB, Leimich P. Approaches to the space-time modelling of infectious disease behaviour. *Math Med Biol* 2000;17(10):1–13.
- Lunn D, Spiegelhalter D, Thomas A, Best N. The BUGS Project: evolution, critique and future directions. *Stat Med* 2009. doi:10.1002/sim.3680.
- Neal P, Roberts G. Statistical inference and model selection for the 1861 Hagelloch measles epidemic. *Biostatistics* 2004;5(2):249–61.
- Morton AM, Finkenstädt BF. Discrete-time modelling of disease incidence time series by using Markov Chain Monte Carlo methods. *Appl Stat* 2005;54(3):575–94.
- Spiegelhalter D, Best N, Carlin B, van der Linde A. Bayesian deviance, the effective number of parameters and the comparison of arbitrarily complex models. *J R Stat Soc B* 2002;64:583–640.