

*Epidemiology and Infection*  
Supplementary material

Quantifying differences in the epidemic curves from three influenza surveillance systems: a nonlinear regression analysis

E.G. Thomas<sup>1†</sup>, J.M. McCaw<sup>1,2</sup>, H.A. Kelly<sup>3,4</sup>, K.A. Grant<sup>3</sup>, J. McVernon<sup>1,2</sup>

† Corresponding author; emma.thomas@unimelb.edu.au

1. Centre for Epidemiology and Biostatistics, Melbourne School of Population and Global Health, University of Melbourne, Victoria, Australia
2. Vaccine and Immunisation Research Group, Murdoch Childrens' Research Institute, Victoria, Australia
3. Epidemiology Unit, Victorian Infectious Diseases Reference Laboratory, Victoria, Australia
4. National Centre for Epidemiology and Population Health, Australian National University, ACT, Australia

## Methodological details: nonlinear regression

Here, we describe in greater detail the nonlinear regression model used to quantify the contribution of surveillance system, year, age group and region to the shape of the influenza epidemic curves. We had  $N = 8,416$  response observations  $y_1, \dots, y_N$  representing counts of incident influenza cases in each surveillance week by year, surveillance system, age group and region. These attributes were encoded as categorical variables in the corresponding covariate observations  $\mathbf{x}_1, \dots, \mathbf{x}_N$ .

We modelled the conditional mean number of incident cases in week  $t$ ,  $\lambda_t(\mathbf{x}, \beta)$ , as a non-linear function  $\lambda_t(\mathbf{x}, \beta) = f(t, \mathbf{x}, \beta)$  where  $Y_t$  is the disease count in week  $t$ ,  $\mathbf{x}$  is a vector of covariate values,  $\beta$  is a parameter vector and  $\mathbf{x}$  is the covariate vector. We specified the function  $f$  to be the sum of a Gaussian function representing the epidemic component and a time-invariant constant representing the endemic component of weekly disease incidence, so that

$$f(t, \mathbf{x}, \beta) = \frac{S(\mathbf{x}, \beta)}{\sqrt{2\pi}D(\mathbf{x}, \beta)/3.92} \exp \left[ -\frac{(t - T(\mathbf{x}, \beta))^2}{2D(\mathbf{x}, \beta)^2} \right] + C(\mathbf{x}, \beta)$$

$S$  may be interpreted as determining the final size of the epidemic component,  $D$  the length of the epidemic component,  $T$  the timing of the week of peak incidence and  $C$  the constant baseline or endemic activity in the system which is present year-round (both in and out of influenza season). As  $D/3.92$  is the standard deviation parameter of the Gaussian function in equation (1) and  $T$  is its mean, we would expect approximately 95% of epidemic cases to occur in the time interval  $T \pm 0.5D$ , so that  $D$  approximates the duration of the seasonal epidemic; we refer to this quantity as the epidemic duration.

Setting  $\beta^T = [\beta_S^T \quad \beta_D^T \quad \beta_T^T \quad \beta_C^T \quad \beta_\alpha^T]$ , we then modelled  $\log S$ ,  $D$ ,  $T$  and  $\log C$  as linear in  $\mathbf{x}$  (where  $\log a$  is the natural logarithm of  $a$ ). We chose to model  $\log S$  and  $\log C$  rather than  $S$  and  $C$  as these parameters represent mean counts of disease and so are necessarily positive. Thus we defined four linear predictors:

$$\log S(\mathbf{x}, \beta) = \mathbf{x}^T \beta_S; \quad D(\mathbf{x}, \beta) = \mathbf{x}^T \beta_D; \quad T(\mathbf{x}, \beta) = \mathbf{x}^T \beta_T; \quad \log C(\mathbf{x}, \beta) = \mathbf{x}^T \beta_C \quad (1)$$

For each linear predictor we also fit interaction terms between all variables and ‘pandemic year’ (2009) as an indicator variable. This decision was made a priori as we judged that the relative effects of surveillance system, age group and region were likely to be very different in a year in which a pandemic strain first appeared as compared to subsequent years, and evidence from the 2009 pandemic supports this.

We expected unmeasured sources of variation in the data to cause over-dispersion in the counts of incident disease. We therefore modelled the conditional weekly counts as independent negative binomial random variables  $(Y_t | \mathbf{x}, \beta) \sim \text{NB}(r(\mathbf{x}, \beta), p_t(\mathbf{x}, \beta))$ , where  $p_t(\mathbf{x}, \beta) = \lambda_t(\mathbf{x}, \beta) / (r(\mathbf{x}, \beta) + \lambda_t(\mathbf{x}, \beta))$  and  $\alpha(\mathbf{x}, \beta) = 1/r(\mathbf{x}, \beta)$  models the over-dispersion (larger values of  $\alpha$  correspond to greater over-dispersion). We expected greater numbers of false positives in the ILI data, leading to additional sources of variability in case counts and thus greater over-dispersion of the ILI compared to the

laboratory confirmed data. We saw no a priori reason to expect over-dispersion to vary according to the other variables in our model. We therefore allowed  $\alpha(\mathbf{x}, \beta) = \mathbf{x}^T \beta_\alpha$  to vary by surveillance system.

For observations  $y_1, \dots, y_N$  and covariates  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , the log-likelihood for our model is given by

$$L = \sum_{i=1}^N r(\mathbf{x}_i, \beta) \log [1 - p_t(\mathbf{x}_i, \beta)] + y_i \log [p_t(\mathbf{x}_i, \beta)] + \log \Gamma \left[ \frac{y_i + r(\mathbf{x}_i, \beta)}{(y_i + r(\mathbf{x}_i, \beta))r(\mathbf{x}_i, \beta)} \right]$$

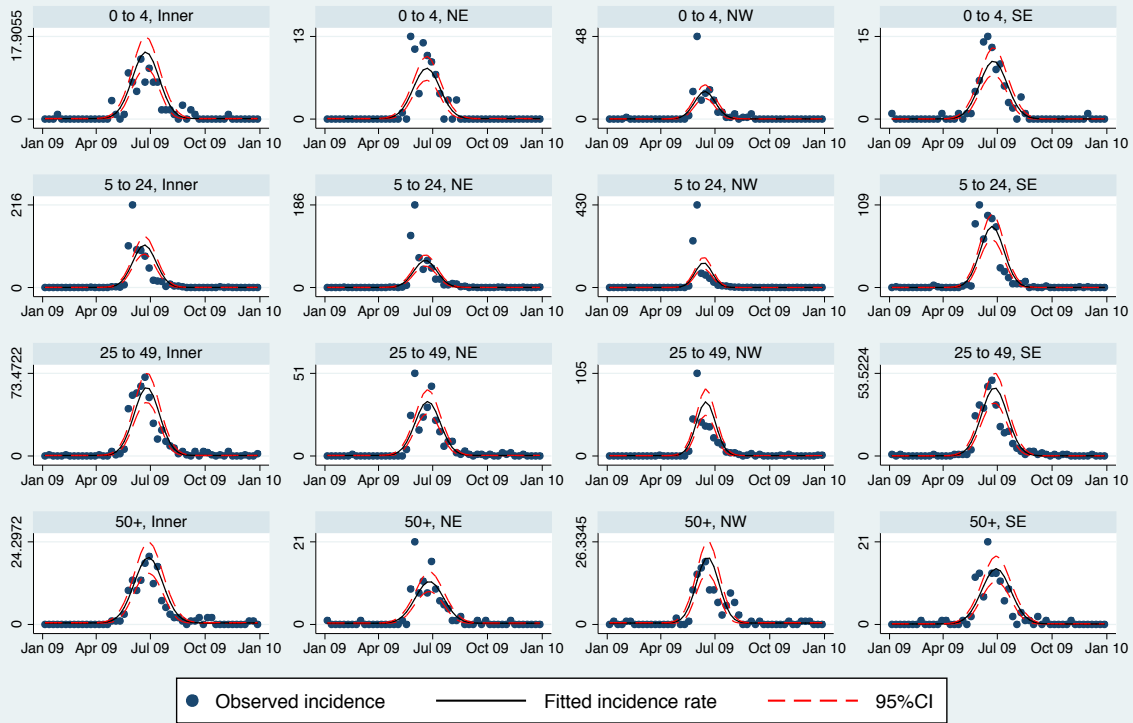
We numerically maximised  $L$  and computed likelihood-based confidence intervals for elements of  $\beta$  using the maximum likelihood estimation command `ml` in Stata.

## Data and nonlinear fits

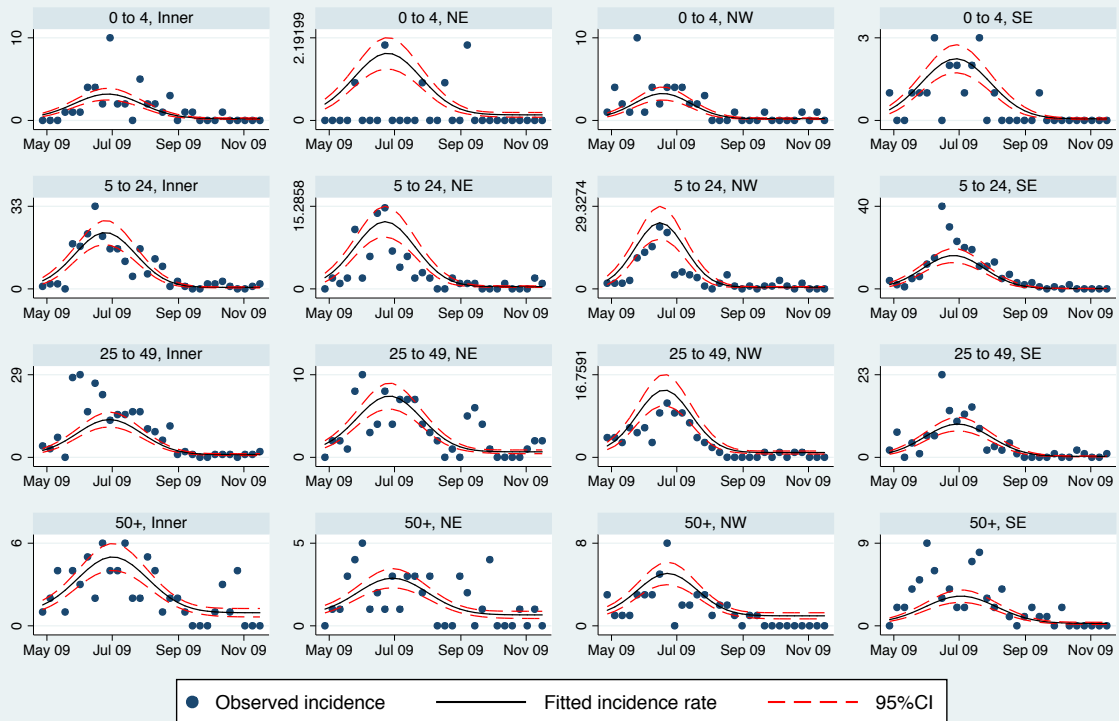
The following plots show the data and resulting fits using the nonlinear model described above. Each plot shows a time series of observed case counts for a particular year (2009–2012) surveillance system (Victorian Department of Health/DH, General Practice Sentinel Surveillance Scheme/GPSS, or Melbourne Medical Deputising Service/MMDS), age group (0–4, 5–24, 25–49 or 50+ years) and region of Melbourne (Inner, North West/NW, North East/NE or South East/SE). The nonlinear model for the mean disease count described by equation (1) and 95%CI is also shown.

The graphs are organised by year and surveillance system. The relevant age group and region are shown above each subplot.

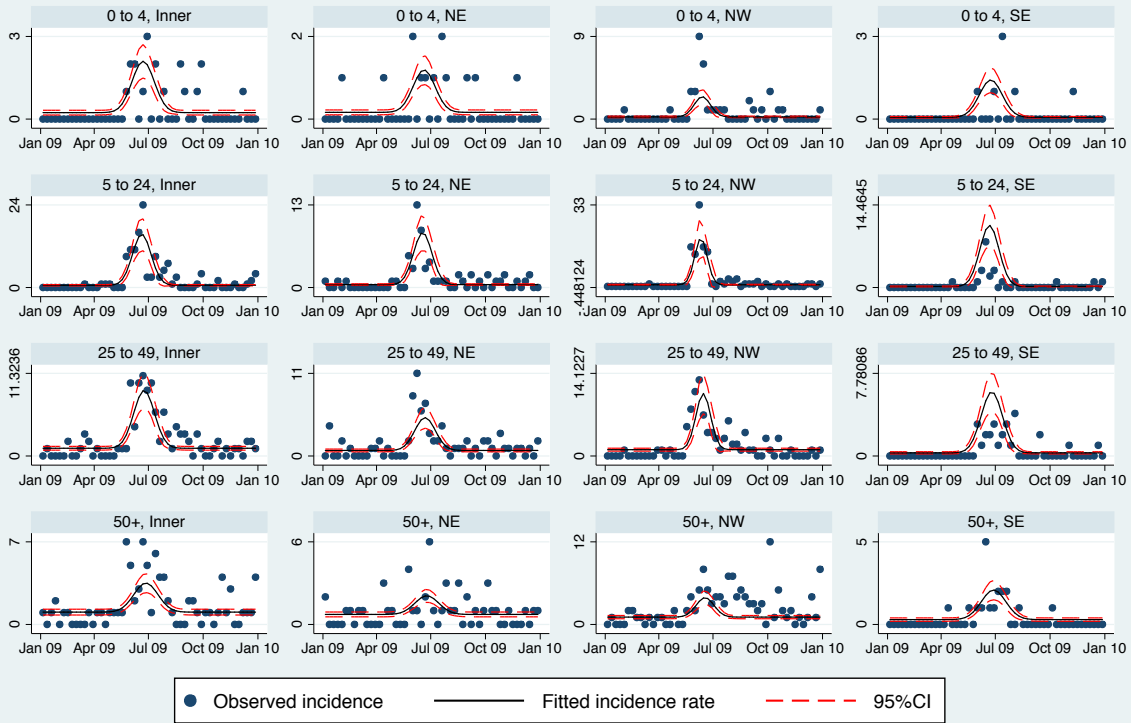
## DH 09



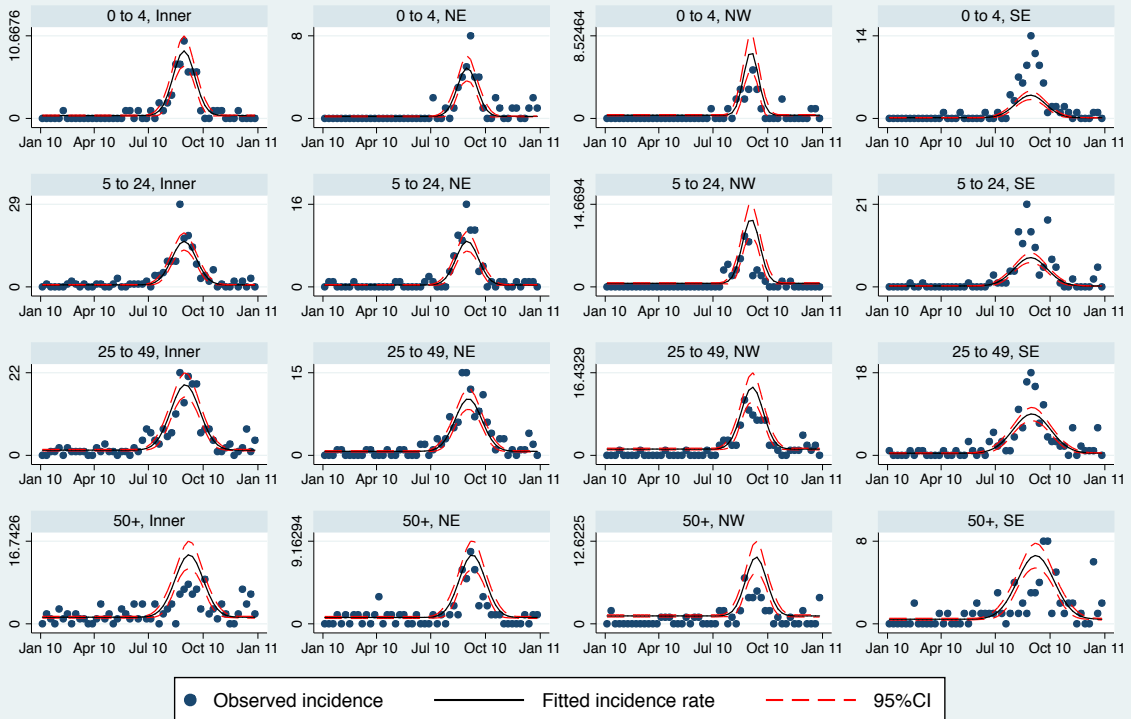
## GPSS 09



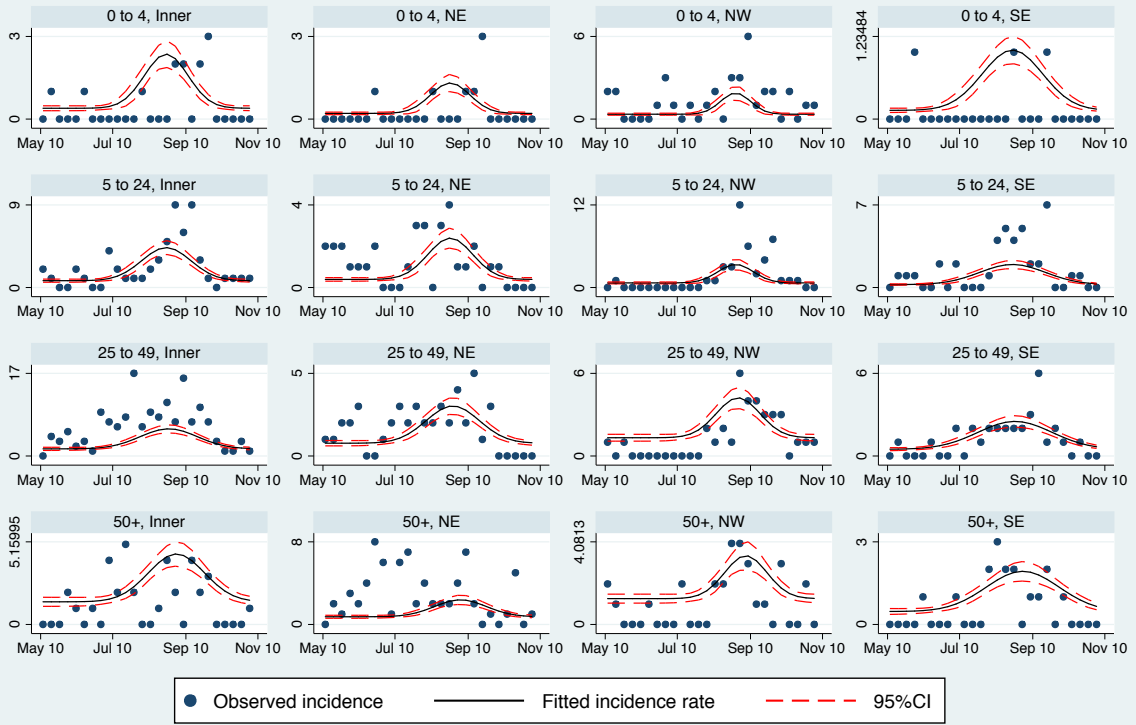
## MMDS 09



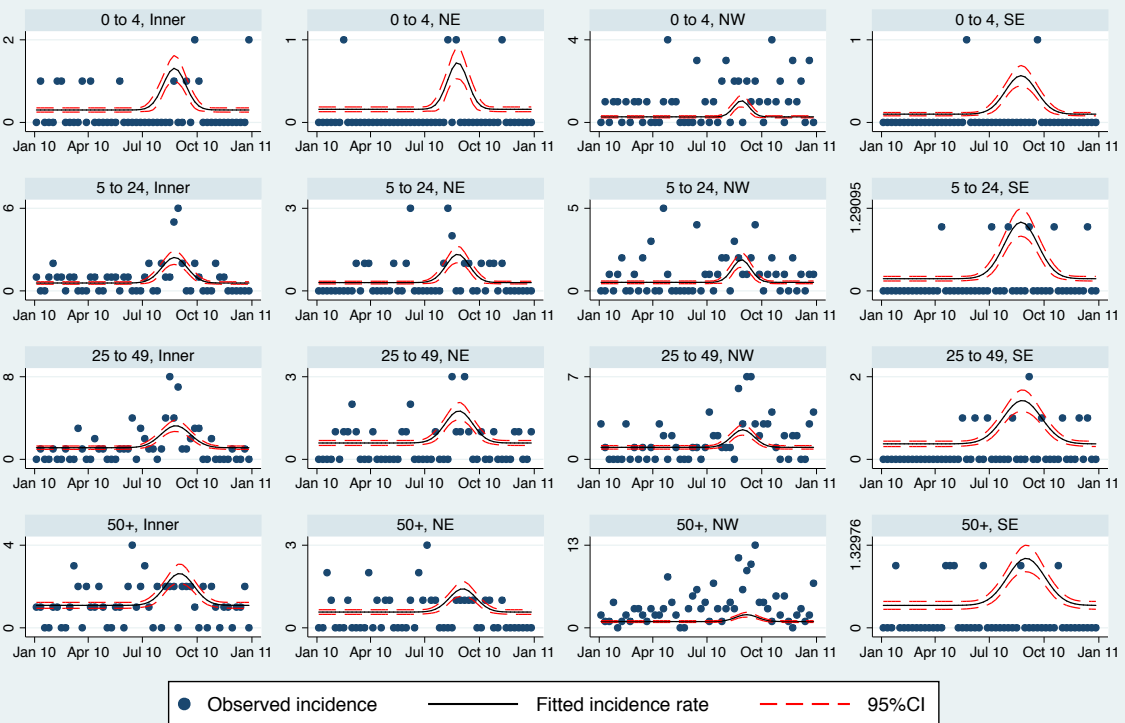
## DH 10



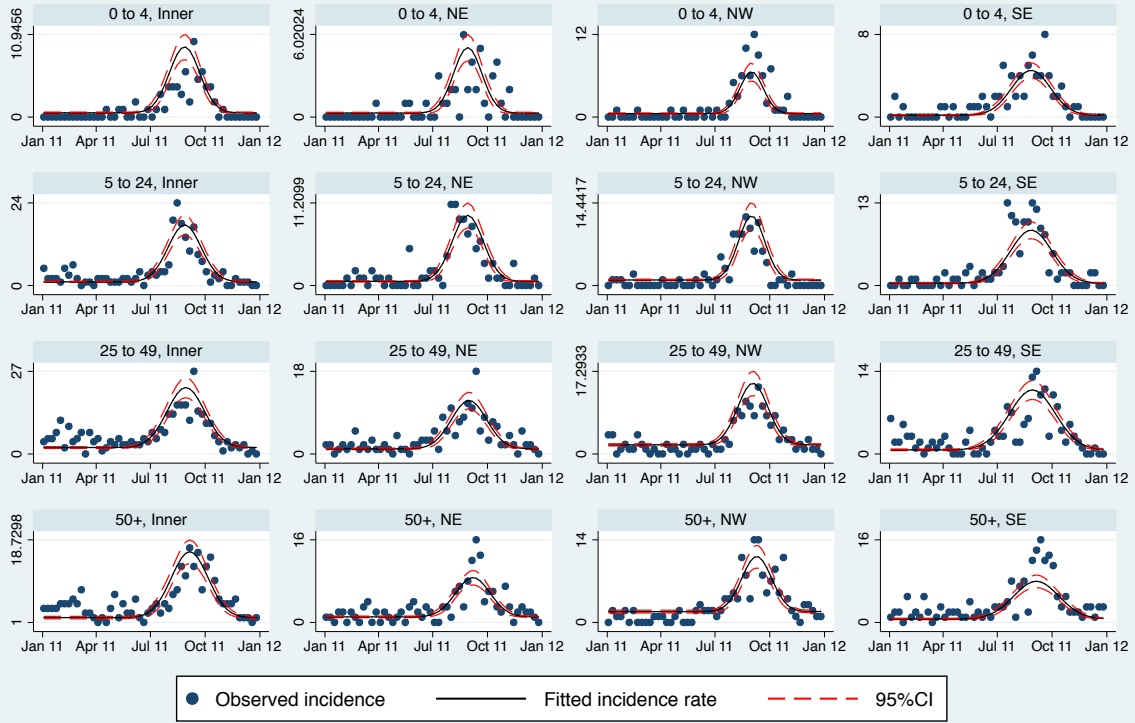
## GPSS 10



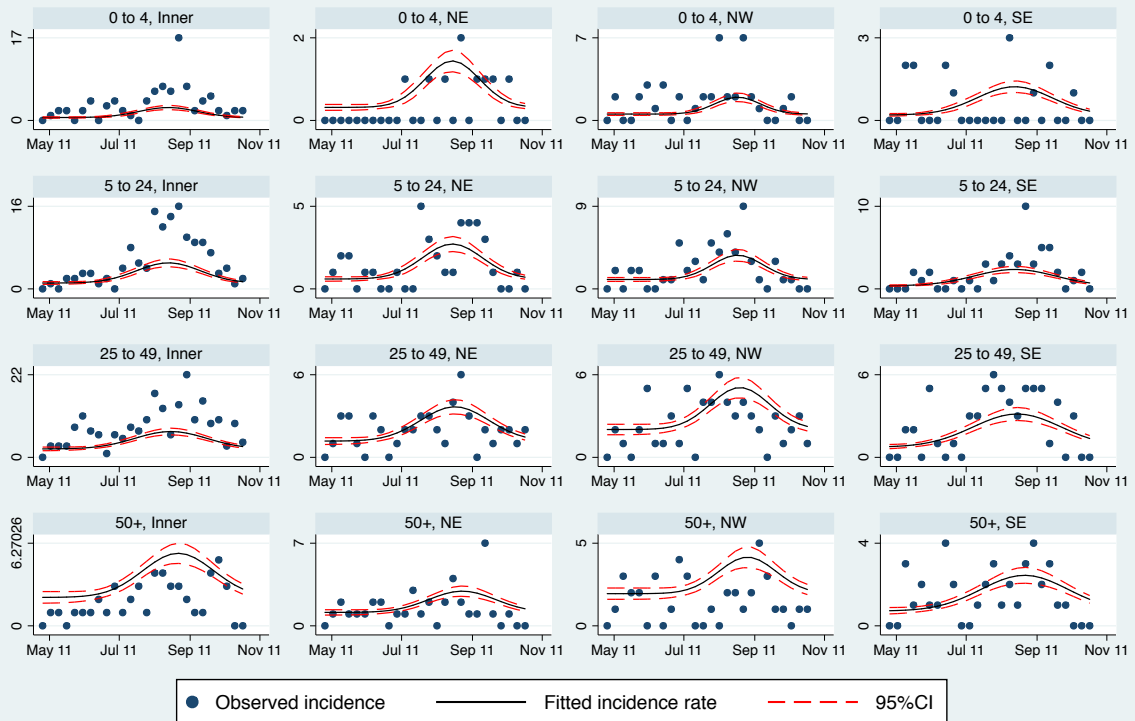
## MMDS 10



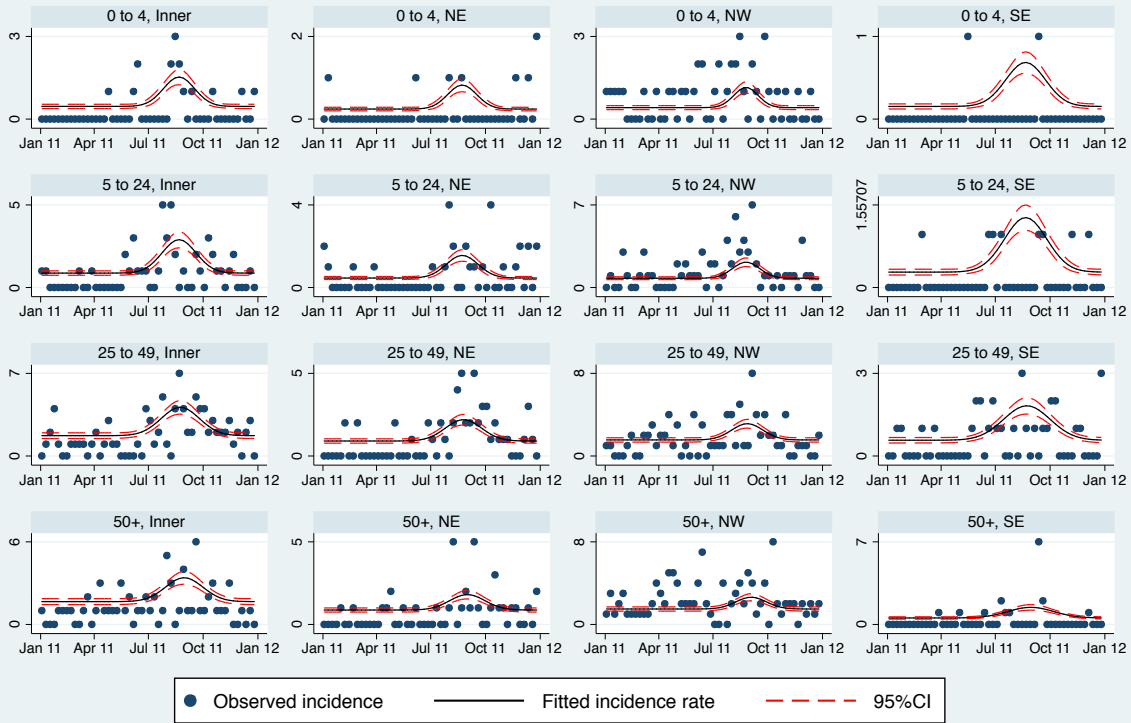
## DH 11



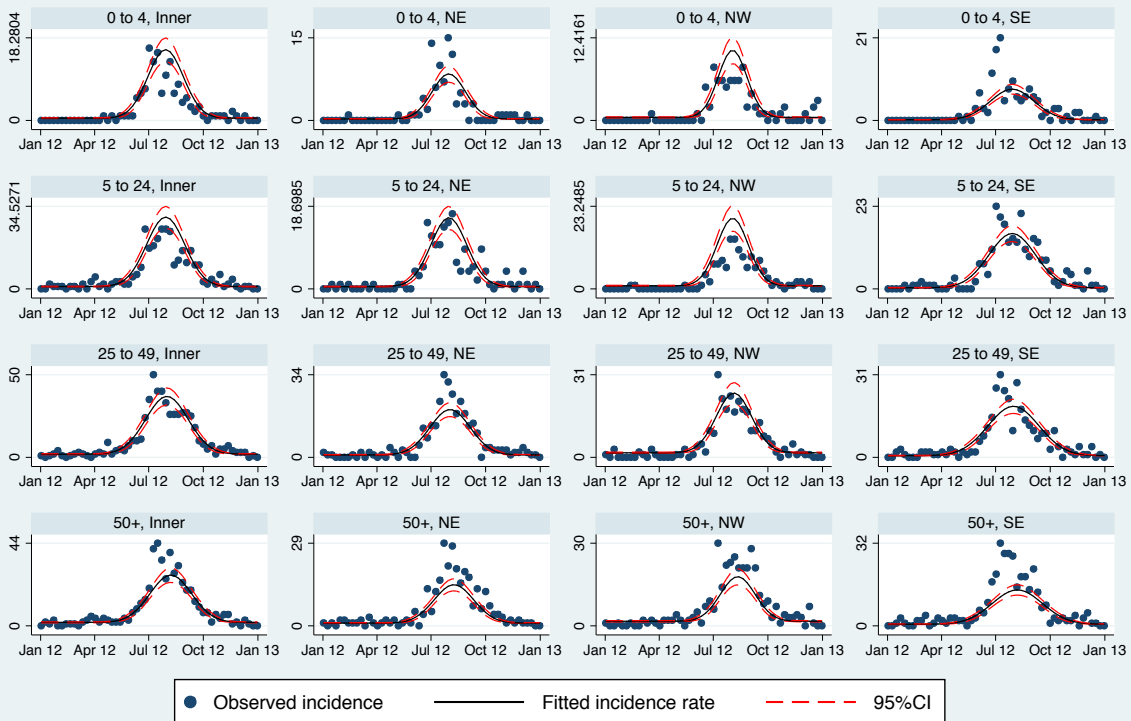
## GPSS 11



## MMDS 11

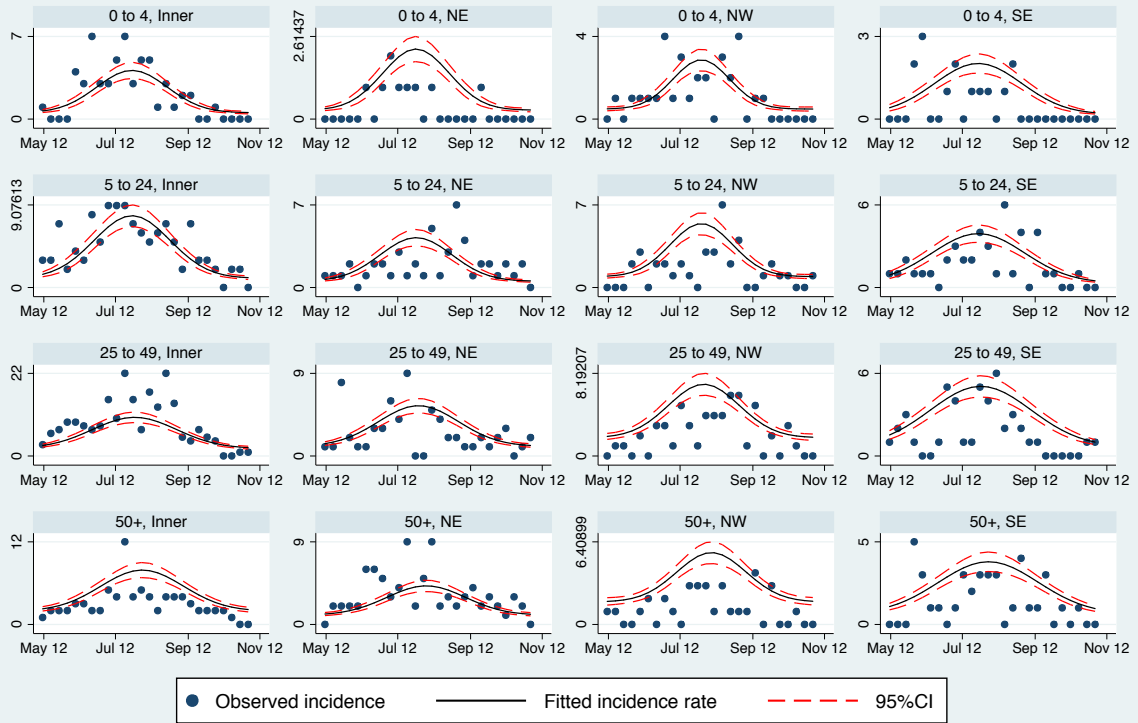


## DH 12





## GPSS 12



## MMDS 12

