

# Bayesian Spatio-Temporal Modelling of Disease Incidence with Nonignorable Missingness

Miles Moran

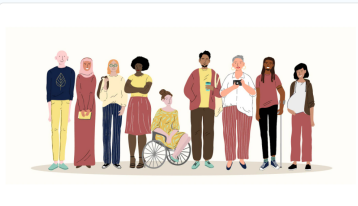
Oregon State University

November 4, 2024

# The Motivating Principle

## CDC's CORE Commitment to Health Equity

CDC works to ensure health equity is embedded in an all-of-public health approach.

[Learn More >](#)

# The Motivating Principle

## CDC's CORE Commitment to Health Equity

CDC works to ensure health equity is embedded in an all-of-public health approach.

[Learn More >](#)



CDC's CORE Commitment to Health Equity stands on four key pillars:

- **Cultivating** comprehensive health equity science
- **Optimizing** interventions
- **Reinforcing** and expanding robust partnerships
- **Enhancing** capacity and workplace diversity, inclusion, and engagement.

# The Motivating Principle

## CDC's CORE Commitment to Health Equity

CDC works to ensure health equity is embedded in an all-of-public health approach.

[Learn More >](#)



CDC's CORE Commitment to Health Equity stands on four key pillars:

- **Cultivating comprehensive health equity science**
- **Optimizing interventions**
- **Reinforcing and expanding robust partnerships**
- **Enhancing capacity and workplace diversity, inclusion, and engagement.**

# What it Takes to Achieve Health Equity

- *More than ever before,*

# What it Takes to Achieve Health Equity

- *More than ever before,*
  - Disease surveillance data are more prevalent (and detailed!)

# What it Takes to Achieve Health Equity

- *More than ever before,*
  - Disease surveillance data are more prevalent (and detailed!)
  - Models of disease spread (e.g. SIR-like) are well-studied

# What it Takes to Achieve Health Equity

- *More than ever before,*
  - Disease surveillance data are more prevalent (and detailed!)
  - Models of disease spread (e.g. SIR-like) are well-studied
  - Modern computational capacity can meet our demands



# What it Takes to Achieve Health Equity

- *More than ever before,*
  - Disease surveillance data are more prevalent (and detailed!)
  - Models of disease spread (e.g. SIR-like) are well-studied
  - Modern computational capacity can meet our demands
- *However,*

# What it Takes to Achieve Health Equity

- *More than ever before,*
  - Disease surveillance data are more prevalent (and detailed!)
  - Models of disease spread (e.g. SIR-like) are well-studied
  - Modern computational capacity can meet our demands
- *However,*
  - The data are *incomplete* for some demographic covariates

# What it Takes to Achieve Health Equity

- *More than ever before,*
  - Disease surveillance data are more prevalent (and detailed!)
  - Models of disease spread (e.g. SIR-like) are well-studied
  - Modern computational capacity can meet our demands
- *However,*
  - The data are *incomplete* for some demographic covariates
  - *If* these covariates are Missing Not-at-Random (MNAR), then models fit with imputed data (or complete-cases-only) will yield biased parameter estimates

# What it Takes to Achieve Health Equity

- *More than ever before,*
  - Disease surveillance data are more prevalent (and detailed!)
  - Models of disease spread (e.g. SIR-like) are well-studied
  - Modern computational capacity can meet our demands
- *However,*
  - The data are *incomplete* for some demographic covariates
  - *If* these covariates are Missing Not-at-Random (MNAR), then models fit with imputed data (or complete-cases-only) will yield biased parameter estimates

## Definition: MNAR Data

With  $p = \Pr(\text{race is observed})$ , we say that the race variable is Missing Not-at-Random (MNAR) if  $p$  depends on race itself or on other *unobserved* variables

# What it Takes to Achieve Health Equity

- *More than ever before,*
  - Disease surveillance data are more prevalent (and detailed!)
  - Models of disease spread (e.g. SIR-like) are well-studied
  - Modern computational capacity can meet our demands
- *However,*
  - The data are *incomplete* for some demographic covariates
  - *If* these covariates are Missing Not-at-Random (MNAR), then models fit with imputed data (or complete-cases-only) will yield biased parameter estimates
  - For many vulnerable or marginalized subpopulations, we can *expect* demographic membership to impact missingness ( $\therefore$  MNAR)

# Our Research Goals

- Develop a spatio-temporal model for the joint distribution of

$X_{tgi j}$  = number of cases from  $(t, g, i)$  observed with race  $j$

$M_{tgi}$  = number of cases from  $(t, g, i)$  missing race

for every time period  $t$ , geographic unit  $g$ , stratum  $i$ , and race  $j$ .

# Our Research Goals

- Develop a spatio-temporal model for the joint distribution of

$X_{tgi j}$  = number of cases from  $(t, g, i)$  observed with race  $j$

$M_{tgi}$  = number of cases from  $(t, g, i)$  missing race

for every time period  $t$ , geographic unit  $g$ , stratum  $i$ , and race  $j$ .

- Conduct a simulation study to demonstrate the model's advantage over SIR-like models that account for missingness in other ways (e.g. imputation or complete-case analysis)

# Our Research Goals

- Develop a spatio-temporal model for the joint distribution of

$X_{tgi j}$  = number of cases from  $(t, g, i)$  observed with race  $j$

$M_{tgi}$  = number of cases from  $(t, g, i)$  missing race

for every time period  $t$ , geographic unit  $g$ , stratum  $i$ , and race  $j$ .

- Conduct a simulation study to demonstrate the model's advantage over SIR-like models that account for missingness in other ways (e.g. imputation or complete-case analysis)
- Apply the model to real-world data (in our case, COVID-19 incidence in Michigan, disaggregated by race/ethnicity, age group, and sex)



# Our Research Goals

- Develop a spatio-temporal model for the joint distribution of

$X_{tgi j}$  = number of cases from  $(t, g, i)$  observed with race  $j$

$M_{tgi}$  = number of cases from  $(t, g, i)$  missing race

for every time period  $t$ , geographic unit  $g$ , stratum  $i$ , and race  $j$ .

- Conduct a simulation study to demonstrate the model's advantage over SIR-like models that account for missingness in other ways (e.g. imputation or complete-case analysis)
- Apply the model to real-world data (in our case, COVID-19 incidence in Michigan, disaggregated by race/ethnicity, age group, and sex)

# Overview: Disease Process Models

- Compartment models where transition rates are based on assumptions of homogeneous mixing and the *law of mass action*

# Overview: Disease Process Models

- Compartment models where transition rates are based on assumptions of homogeneous mixing and the *law of mass action*
- Deterministic Models: ODE-based, cts. in time. With infection rate  $\beta$  and recovery rate  $\gamma$ , under frequency-dependent transmission:

$$\begin{aligned}\frac{dX(t)}{dt} &= -\frac{\beta X(t)Y(t)}{N} \\ \frac{dY(t)}{dt} &= \frac{\beta X(t)Y(t)}{N} - \gamma Y(t) \\ \frac{dZ(t)}{dt} &= \gamma Y(t)\end{aligned}$$

# Overview: Disease Process Models

- Compartment models where transition rates are based on assumptions of homogeneous mixing and the *law of mass action*
- Deterministic Models: ODE-based, cts. in time. With infection rate  $\beta$  and recovery rate  $\gamma$ , under frequency-dependent transmission:

$$\begin{aligned}\frac{dX(t)}{dt} &= -\frac{\beta X(t)Y(t)}{N} \\ \frac{dY(t)}{dt} &= \frac{\beta X(t)Y(t)}{N} - \gamma Y(t) \\ \frac{dZ(t)}{dt} &= \gamma Y(t)\end{aligned}$$

- Working with surveillance (i.e. population-level) data, We want a *stochastic, discrete-time* analog to this SIR model, e.g.

$$\begin{aligned}X_t &= X_{t-1} + B_{t-d} - Y_t \\ Y_t &\sim F(\dots)\end{aligned}$$

# Brief History of Relevant Stochastic Models

## (1) Time-Series SIR ("TSIR") Models:

- Kendall (1949) *Stochastic Processes and Population Growth*.
- Bartlett (1956) *Deterministic and Stochastic Models for Recurrent Epidemics*.
- Bjørnstad, Finkenstädt, Grenfell (2002) *Dynamics of Measles Epidemics: Estimating Scaling of Transmission Rates Using a Time Series SIR Model*

# Brief History of Relevant Stochastic Models

## (1) Time-Series SIR ("TSIR") Models:

- Kendall (1949) *Stochastic Processes and Population Growth*.
- Bartlett (1956) *Deterministic and Stochastic Models for Recurrent Epidemics*.
- Bjørnstad, Finkenstädt, Grenfell (2002) *Dynamics of Measles Epidemics: Estimating Scaling of Transmission Rates Using a Time Series SIR Model*

## (2) Epidemic / Endemic (hhh4/surveillance) Models:

- Held, Höhle, Hofmann (2005) *A Statistical Framework for the Analysis of Multivariate Infectious Disease Surveillance Counts*
- Held & Paul (2012) *Modeling Seasonality in Space-Time Infectious Disease Surveillance Data*

# Brief History of Relevant Stochastic Models

## (1) Time-Series SIR ("TSIR") Models:

- Kendall (1949) *Stochastic Processes and Population Growth*.
- Bartlett (1956) *Deterministic and Stochastic Models for Recurrent Epidemics*.
- Bjørnstad, Finkenstädt, Grenfell (2002) *Dynamics of Measles Epidemics: Estimating Scaling of Transmission Rates Using a Time Series SIR Model*

## (2) Epidemic / Endemic (hhh4/surveillance) Models:

- Held, Höhle, Hofmann (2005) *A Statistical Framework for the Analysis of Multivariate Infectious Disease Surveillance Counts*
- Held & Paul (2012) *Modeling Seasonality in Space-Time Infectious Disease Surveillance Data*

## (3) Zelner Contact-Heterogeneity Models:

- Lloyd-Smith, et al. (2005) *Superspreading and the Effect of Individual Variation on Disease Emergence*
- Zelner et. al (2020) *Understanding the Importance of Contact Heterogeneity and Variable Infectiousness in the Dynamics of a Large Norovirus Outbreak*

# Derivation #1: TSIR Models (1/2)

(1) Let  $Y_t$  denote # infected at time  $t$ . Denote  $Y_0 = y_0$ .



# Derivation #1: TSIR Models (1/2)

- (1) Let  $Y_t$  denote # infected at time  $t$ . Denote  $Y_0 = y_0$ .
- (2) Suppose that each initial infected independently generates a tree of new infections according to a *linear birth process* over  $(0, t)$ :

$$\Pr(\text{new infection in } (t, t + dt)) = \lambda dt + o(dt)$$

# Derivation #1: TSIR Models (1/2)

- (1) Let  $Y_t$  denote # infected at time  $t$ . Denote  $Y_0 = y_0$ .
- (2) Suppose that each initial infected independently generates a tree of new infections according to a *linear birth process* over  $(0, t)$ :

$$\Pr(\text{new infection in } (t, t + dt)) = \lambda dt + o(dt)$$

- (3) Kendall (1949) show that the number of new infections from tree  $i$  is  $T_i$  (where a “success” = the single bernoulli event of *not* infecting someone):

$$(T_1, \dots, T_{y_0}) \stackrel{\text{iid}}{\sim} \text{Geom}\left( \underbrace{e^{-\lambda t}}_{\text{success prob}} \right)$$

# Derivation #1: TSIR Models (1/2)

- (1) Let  $Y_t$  denote # infected at time  $t$ . Denote  $Y_0 = y_0$ .
- (2) Suppose that each initial infected independently generates a tree of new infections according to a *linear birth process* over  $(0, t)$ :

$$\Pr(\text{new infection in } (t, t + dt)) = \lambda dt + o(dt)$$

- (3) Kendall (1949) show that the number of new infections from tree  $i$  is  $T_i$  (where a “success” = the single bernoulli event of *not* infecting someone):

$$(T_1, \dots, T_{y_0}) \stackrel{\text{iid}}{\sim} \text{Geom}\left( \underbrace{e^{-\lambda t}}_{\text{success prob}} \right)$$

- (4) Following the properties of the Negative Binomial distribution,

$$\begin{aligned} \underbrace{(Y_t - y_0)}_{\# \text{ failures}} &= \left( \sum_{i=1}^{y_0} T_i \right) \sim \text{NB}\left( \underbrace{y_0}_{\# \text{ successes}}, \underbrace{e^{-\lambda t}}_{\text{success prob}} \right) \sim \text{NB}\left( \underbrace{y_0(e^{\lambda t} - 1)}_{\text{mean}}, \underbrace{y_0}_{\text{dispersion}} \right) \\ \underbrace{Y_t}_{\# \text{ trials}} &\sim \text{NB}\left( \underbrace{y_0}_{\# \text{ successes}}, \underbrace{e^{-\lambda t}}_{\text{success prob}} \right) \sim \text{NB}\left( \underbrace{y_0 e^{\lambda t}}_{\text{mean}}, \underbrace{y_0}_{\text{dispersion}} \right) \end{aligned}$$

# Derivation #1: TSIR Models (2/2)

- (5) Generalizing this formulation from the interval  $(0, t)$  to our discretized intervals  $(t-1, t)$ , we can write

$$(Y_t - y_{t-1} \mid Y_{t-1} = y_{t-1}) \sim \text{NB}(y_{t-1}(e^\lambda - 1), y_{t-1})$$

and, if we make the simplifying assumption that prevalence = incidence (i.e. all infections recover before the next time period), then just

$$(Y_t \mid Y_{t-1} = y_{t-1}) \sim \text{NB}(y_{t-1}(e^\lambda - 1), y_{t-1})$$

# Derivation #1: TSIR Models (2/2)

- (5) Generalizing this formulation from the interval  $(0, t)$  to our discretized intervals  $(t-1, t)$ , we can write

$$(Y_t - y_{t-1} \mid Y_{t-1} = y_{t-1}) \sim \text{NB}(y_{t-1}(e^\lambda - 1), y_{t-1})$$

and, if we make the simplifying assumption that prevalence = incidence (i.e. all infections recover before the next time period), then just

$$(Y_t \mid Y_{t-1} = y_{t-1}) \sim \text{NB}(y_{t-1}(e^\lambda - 1), y_{t-1})$$

- (6) At this point, there's room for creativity in describing the hazard rate  $\lambda$ . Based on mass-action, we could take, e.g.,  $\lambda = \frac{\beta x_{t-1}}{N}$ . so that

$$\begin{aligned} E(Y_t \mid Y_{t-1} = y_{t-1}, X_{t-1} = x_{t-1}) &= y_{t-1} e^{\frac{\beta x_{t-1}}{N}} - y_{t-1} \\ &\approx \frac{\beta x_{t-1} y_{t-1}}{N} - y_{t-1} \quad \text{for small } \beta x_{t-1}/N \end{aligned}$$

## Derivation #2: Epidemic-Endemic / hhh4 Models

- (1) Assume prevalence = incidence as in TSIR (fixed, unit-length infection time)

## Derivation #2: Epidemic-Endemic / hhh4 Models

- (1) Assume prevalence = incidence as in TSIR (fixed, unit-length infection time)
- (2) Assume basic survival model for susceptibles: constant hazard for each time period so that the time until a susceptible is infected is exponential:

$$\Pr(\text{infection in } (t-1, t] \mid \text{no infection by } t-1) = 1 - e^{-\lambda_t}$$

## Derivation #2: Epidemic-Endemic / hhh4 Models

- (1) Assume prevalence = incidence as in TSIR (fixed, unit-length infection time)
- (2) Assume basic survival model for susceptibles: constant hazard for each time period so that the time until a susceptible is infected is exponential:

$$\Pr(\text{infection in } (t-1, t] \mid \text{no infection by } t-1) = 1 - e^{-\lambda_t}$$

- (3) Treating each susceptible as an independent Bernoulli trial, the number of new cases at time  $t$  becomes

$$(Y_t \mid X_{t-1} = x_{t-1}) \sim \text{Binom}(x_{t-1}, 1 - e^{-\lambda_t})$$
$$\dot{\sim} \text{Pois}(x_{t-1} \lambda_t) \quad \text{when } x_{t-1} \text{ large and } \lambda_t \text{ small}$$



# Derivation #2: Epidemic-Endemic / hhh4 Models

- (1) Assume prevalence = incidence as in TSIR (fixed, unit-length infection time)
- (2) Assume basic survival model for susceptibles: constant hazard for each time period so that the time until a susceptible is infected is exponential:

$$\Pr(\text{infection in } (t-1, t] \mid \text{no infection by } t-1) = 1 - e^{-\lambda_t}$$

- (3) Treating each susceptible as an independent Bernoulli trial, the number of new cases at time  $t$  becomes

$$(Y_t \mid X_{t-1} = x_{t-1}) \sim \text{Binom}(x_{t-1}, 1 - e^{-\lambda_t})$$

$$\dot{\sim} \text{Pois}(x_{t-1} \lambda_t) \quad \text{when } x_{t-1} \text{ large and } \lambda_t \text{ small}$$

- (4) Again, there's room for creativity in describing the hazard rate  $\lambda$ . Based on mass-action, we arrive at

$$(Y_t \mid Y_{t-1} = y_{t-1}, X_{t-1} = x_{t-1},) \dot{\sim} \text{Pois}\left(\frac{\beta x_{t-1} y_{t-1}}{N}\right)$$

# Our Research Goals

- Develop a spatio-temporal model for the joint distribution of

$X_{tgi j}$  = number of cases from  $(t, g, i)$  observed with race  $j$

$M_{tgi}$  = number of cases from  $(t, g, i)$  missing race

for every time period  $t$ , geographic unit  $g$ , stratum  $i$ , and race  $j$ .

- Conduct a simulation study to demonstrate the model's advantage over SIR-like models that account for missingness in other ways (e.g. imputation or complete-case analysis)
- Apply the model to real-world data (in our case, COVID-19 incidence in Michigan, disaggregated by race/ethnicity, age group, and sex)

# Missingness in the Simplest Case

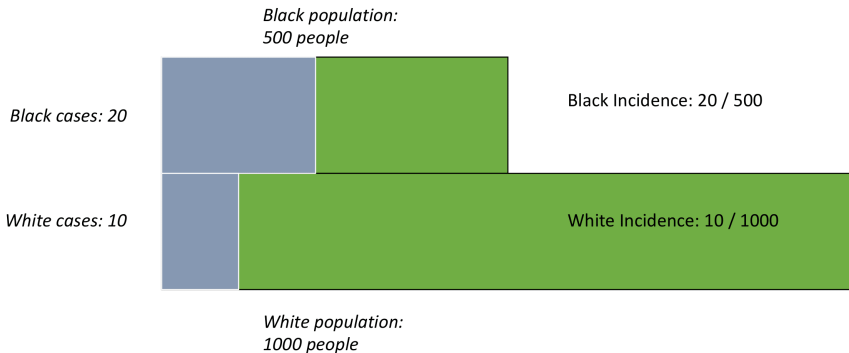
*Black population:  
500 people*



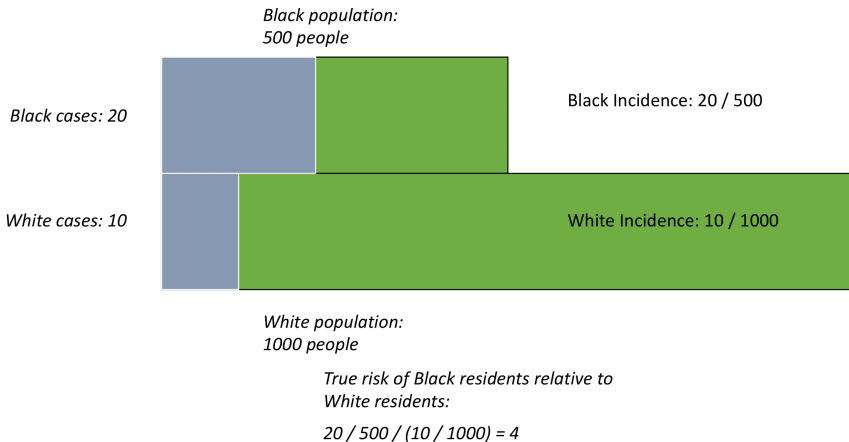
*White population:  
1000 people*



# Missingness in the Simplest Case



# Missingness in the Simplest Case



# Missingness in the Simplest Case

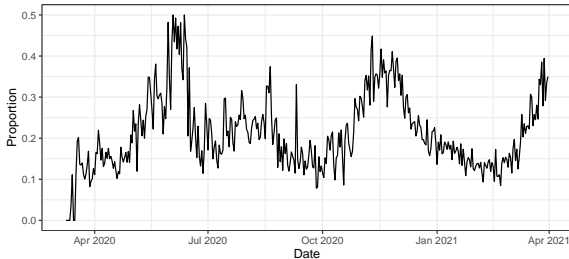


# Missingness in the Simplest Case

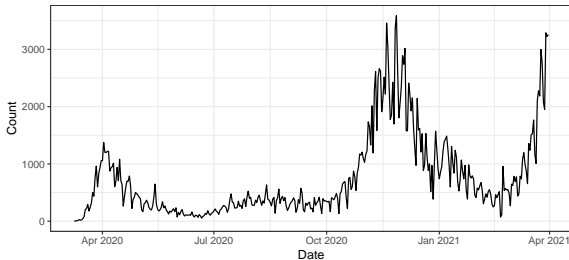


# Missingness in Real Life

Proportion cases missing race by date



Number incident cases by date





# Overview: Missingness Process Models

- At time  $t$ , for individual  $k$  from race  $j$ , let
  - $Y_{tjk}$  denote presence/absence of infection (binary)
  - $M_{tjk}$  denote missingness (binary)so that  $Y_{tjk}$  is observed only when  $M_{tjk} = 1$
- A model for data subject to missingness is just a specification of  $f(\mathbf{y}, \mathbf{m})$ , which is done through some kind of decomposition:

Selection Factorization:  $f(\mathbf{y}, \mathbf{m}) = f_1(\mathbf{y})f_2(\mathbf{m}|\mathbf{y})$

Pattern-Mixture Factorization:  $f(\mathbf{y}, \mathbf{m}) = g_1(\mathbf{m})g_2(\mathbf{y}|\mathbf{m})$

Random Effects:  $f(\mathbf{y}, \mathbf{m}) = \int h_1(\mathbf{y}|\mathbf{u})h_2(\mathbf{m}|\mathbf{u})h_3(\mathbf{u})d\mathbf{u}$

## Derivation #3: Missingness in a TSIR-like Model

- (1) Model *true* incidence-by-race with a discrete-time 1st-order Markov model, i.e.  $(Y_{tgi j} | \mathbf{Y}_{(t-1)}) \sim \text{Pois}(\lambda_{tgi j}^{\text{TOT}})$

## Derivation #3: Missingness in a TSIR-like Model

- (1) Model *true* incidence-by-race with a discrete-time 1st-order Markov model, i.e.  $(Y_{tgi j} | \mathbf{Y}_{(t-1)}) \sim \text{Pois}(\lambda_{tgi j}^{\text{TOT}})$
- (2) Model *observed* incidence-by-race as independent draws from the population of infected, i.e.  $(X_{tgi j} | Y_{tgi j}) \sim \text{Binom}(Y_{tgi j}, p_{tgi j})$

## Derivation #3: Missingness in a TSIR-like Model

- (1) Model *true* incidence-by-race with a discrete-time 1st-order Markov model, i.e.  $(Y_{tgi j} | \mathbf{Y}_{(t-1)}) \sim \text{Pois}(\lambda_{tgi j}^{\text{TOT}})$
- (2) Model *observed* incidence-by-race as independent draws from the population of infected, i.e.  $(X_{tgi j} | Y_{tgi j}) \sim \text{Binom}(Y_{tgi j}, p_{tgi j})$
- (3) Marginalize over  $Y_{tgi j}$  to obtain the *observational* model:

$$\left\{ \begin{array}{l} X_{tgi j} \sim \text{Pois}(\lambda_{tgi j}^{\text{TOT}} p_{tgi j}) \\ M_{tgi} = Y_{tgi \bullet} - \sum_{j=1}^J X_{tgi j} \\ \sim \text{Pois}\left(\sum_{j=1}^J \lambda_{tgi j}^{\text{TOT}} (1 - p_{tgi j})\right) \end{array} \right\}$$

# EE-Like Observational Model Specification

$$(X_{tgi j} | \mathbf{Y}_{(t-1)gi \bullet}) \sim \text{Pois}(\lambda_{tgi j}^{\text{TOT}} p_{tgi j})$$

$$(M_{tgi} | \mathbf{Y}_{(t-1)gi \bullet}) \sim \text{Pois}\left(\sum_{j=1}^J \lambda_{tgi j}^{\text{TOT}} (1 - p_{tgi j})\right)$$

$$\lambda_{tgi j}^{\text{TOT}} = \lambda_{tgi j}^{\text{AR}} Y_{(t-1)gi \bullet} + \lambda_{tgi j}^{\text{NE}} \sum_{g'=1}^G w_{gg'} Y_{(t-1)g' \bullet \bullet} + \lambda_{tgi j}^{\text{EN}} E_{gij}$$

$$\log(\lambda_{tgi j}^{\text{AR}}) = \mu^{\text{AR}} + \alpha_j^{\text{AR}} + \beta_g^{\text{AR}}$$

$$\log(\lambda_{tgi j}^{\text{NE}}) = \mu^{\text{NE}} + \alpha_j^{\text{NE}} + \beta_g^{\text{NE}}$$

$$\begin{aligned} \log(\lambda_{tgi j}^{\text{EN}}) = & \mu^{\text{EN}} + \alpha_j^{\text{EN}} + \beta_g^{\text{EN}} + \gamma^{\text{EN}} t + \delta^{\text{EN}} \sin\left(\frac{t}{52} 2\pi\right) \\ & + \varepsilon^{\text{EN}} \cos\left(\frac{t}{52} 2\pi\right) \end{aligned}$$

$$\text{logit}(p_{tgi j}) = \mu^{(\text{p})} + \alpha_j^{(\text{p})} + \beta_g^{(\text{p})}$$

# EE-Like Observational Model Specification

$$(X_{tgi j} | \mathbf{Y}_{(t-1)gi\bullet}) \sim \text{Pois}(\lambda_{tgi j}^{\text{TOT}} p_{tgi j})$$

$$(M_{tgi} | \mathbf{Y}_{(t-1)gi\bullet}) \sim \text{Pois}\left(\sum_{j=1}^J \lambda_{tgi j}^{\text{TOT}} (1 - p_{tgi j})\right)$$

$$\lambda_{tgi j}^{\text{TOT}} = \lambda_{tgi j}^{\text{AR}} Y_{(t-1)gi\bullet} + \lambda_{tgi j}^{\text{NE}} \sum_{g'=1}^G w_{gg'} Y_{(t-1)g'\bullet\bullet} + \lambda_{tgi j}^{\text{EN}} E_{gij}$$

$$\log(\lambda_{tgi j}^{\text{AR}}) = \mu^{\text{AR}} + \alpha_j^{\text{AR}} + \beta_g^{\text{AR}}$$

$$\log(\lambda_{tgi j}^{\text{NE}}) = \mu^{\text{NE}} + \alpha_i^{\text{NE}} + \beta_a^{\text{NE}}$$

## The Data (Part 1)

$(t, g, i, j) = (\text{time}, \text{location}, \text{stratum}, \text{race})$

$X_{tgi j}$  = number of cases from  $(t, g, i)$  observed with race  $j$

$M_{tgi}$  = number of cases from  $(t, g, i)$  missing race

# EE-Like Observational Model Specification

$$(X_{tgi\bullet} | \mathbf{Y}_{(t-1)gi\bullet}) \sim \text{Pois}(\lambda_{tgij}^{\text{TOT}} p_{tgj})$$

$$(M_{tgi\bullet} | \mathbf{Y}_{(t-1)gi\bullet}) \sim \text{Pois}\left(\sum_{j=1}^J \lambda_{tgij}^{\text{TOT}} (1 - p_{tgj})\right)$$

$$\lambda_{tgij}^{\text{TOT}} = \lambda_{tgj}^{\text{AR}} Y_{(t-1)gi\bullet} + \lambda_{tgj}^{\text{NE}} \sum_{g'=1}^G w_{gg'} Y_{(t-1)g'\bullet\bullet} + \lambda_{tgj}^{\text{EN}} E_{gij}$$

$$\log(\lambda_{tgij}^{\text{AR}}) = \mu^{\text{AR}} + \alpha_i^{\text{AR}} + \beta_a^{\text{AR}}$$

$$\log(\lambda_{tgij}^{\text{NE}}) = \mu^{\text{NE}} + \alpha_i^{\text{NE}} + \beta_a^{\text{NE}}$$

$$\log(\lambda_{tgij}^{\text{EN}}) = \mu^{\text{EN}} + \alpha_i^{\text{EN}} + \beta_a^{\text{EN}}$$

$$\text{logit}(\mathbf{Y}_{tgi\bullet}) = (\text{logit}(Y_{tg1\bullet}), \text{logit}(Y_{tg2\bullet}), \dots, \text{logit}(Y_{tgI\bullet}))$$

**Note:**  $Y_{tgi\bullet}, Y_{tg\bullet\bullet}$  are observed, but  $Y_{tgij}$  aren't

# EE-Like Observational Model Specification

$$(X_{tgi j} | \mathbf{Y}_{(t-1)gi\bullet}) \sim \text{Pois}(\lambda_{tgi j}^{\text{TOT}} p_{tgi j})$$

$$(M_{tgi} | \mathbf{Y}_{(t-1)gi\bullet}) \sim \text{Pois}\left(\sum_{j=1}^J \lambda_{tgi j}^{\text{TOT}} (1 - p_{tgi j})\right)$$

$$\lambda_{tgi j}^{\text{TOT}} = \lambda_{tgi j}^{\text{AR}} Y_{(t-1)gi\bullet} + \lambda_{tgi j}^{\text{NE}} \sum_{g'=1}^G w_{gg'} Y_{(t-1)g'\bullet\bullet} + \lambda_{tgi j}^{\text{EN}} E_{gij}$$

$$\log(\lambda_{tgi j}^{\text{AR}}) = \mu^{\text{AR}} + \alpha_i^{\text{AR}} + \beta^{\text{AR}}$$

$$\log(\lambda_{tgi j}^{\text{NE}})$$

$$\log(\lambda_{tgi j}^{\text{EN}})$$

$$\text{logit}(p_{tgi j})$$

## The Data (Part 3)

$w_{gg'}$  = distance weight between locations  $g$  and  $g'$

$$= \begin{cases} \frac{1}{(\# \text{ neighbors})_g} & \text{if } g \text{ adjacent to } g' \\ 0 & \text{if } g \text{ not adjacent to } g' \text{ (or } g = g') \end{cases}$$

$E_{gij}$  = population count from  $(g, i, j)$



# EE-Like Observational Model Specification

$$(X_{tgi j} | \mathbf{Y}_{(t-1)gi \bullet}) \sim \text{Pois}(\lambda_{tgi j}^{\text{TOT}} p_{tgi j})$$

$$(M_{tgi} | \mathbf{Y}_{(t-1)gi \bullet}) \sim \text{Pois}\left(\sum_{j=1}^J \lambda_{tgi j}^{\text{TOT}} (1 - p_{tgi j})\right)$$

$$\lambda_{tgi j}^{\text{TOT}} = \lambda_{tgi j}^{\text{AR}} Y_{(t-1)gi \bullet} + \lambda_{tgi j}^{\text{NE}} \sum_{g'=1}^G w_{gg'} Y_{(t-1)g' \bullet \bullet} + \lambda_{tgi j}^{\text{EN}} E_{gi j}$$

$$\log(\lambda_{tgi j}^{\text{AR}}) = \mu^{\text{AR}} + \alpha_i^{\text{AR}} + \beta_a^{\text{AR}}$$

The Disease Process (Competing Risks Framework)

$\lambda_{tgi j}^{\text{TOT}}$  = hazard rate from all sources

$\lambda_{tgi j}^{\text{AR}}$  = hazard rate of self-area/“autoregressive” infections

$\lambda_{tgi j}^{\text{NE}}$  = hazard rate of neighboring-area infections

$\lambda_{tgi j}^{\text{EN}}$  = hazard rate of background/“environmental” infection

# EE-Like Observational Model Specification

$$(X_{tgi j} | \mathbf{Y}_{(t-1)gi \bullet}) \sim \text{Pois}(\lambda_{tgi j}^{\text{TOT}} p_{tgi j})$$

$$(M_{tgi} | \mathbf{Y}_{(t-1)gi \bullet}) \sim \text{Pois}\left(\sum_{j=1}^J \lambda_{tgi j}^{\text{TOT}} (1 - p_{tgi j})\right)$$

$$\lambda_{tgi j}^{\text{TOT}} = \lambda_{tgi j}^{\text{AR}} Y_{(t-1)gi \bullet} + \lambda_{tgi j}^{\text{NE}} \sum_{g'=1}^G w_{gg'} Y_{(t-1)g' \bullet \bullet} + \lambda_{tgi j}^{\text{EN}} E_{gij}$$

$$\log(\lambda_{tgi j}^{\text{AR}}) = \mu^{\text{AR}} + \alpha_j^{\text{AR}} + \beta_g^{\text{AR}}$$

$$\log(\lambda_{tgi j}^{\text{NE}}) = \mu^{\text{NE}} + \alpha_j^{\text{NE}} + \beta_g^{\text{NE}}$$

$$\log(\lambda_{tgi j}^{\text{EN}}) = \mu^{\text{EN}} + \alpha_j^{\text{EN}} + \beta_g^{\text{EN}} + \gamma^{\text{EN}} t + \delta^{\text{EN}} \sin\left(\frac{t}{52} 2\pi\right) + \varepsilon^{\text{EN}} \cos\left(\frac{t}{52} 2\pi\right)$$

## The Missingness Process

$p_{tgi j}$  = probability that a case from  $(t, g, j)$  reports their race

# EE-Like Observational Model Specification

$$(X_{tgi j} | \mathbf{Y}_{(t-1)gi \bullet}) \sim \text{Pois}(\lambda_{tgi j}^{\text{TOT}} p_{tgi j})$$

$$(M_{tgi} | \mathbf{Y}_{(t-1)gi \bullet}) \sim \text{Pois}\left(\sum_{j=1}^J \lambda_{tgi j}^{\text{TOT}} (1 - p_{tgi j})\right)$$

$$\lambda_{tgi j}^{\text{TOT}} = \lambda_{tgi j}^{\text{AR}} Y_{(t-1)gi \bullet} + \lambda_{tgi j}^{\text{NE}} \sum_{g'=1}^G w_{gg'} Y_{(t-1)g' \bullet \bullet} + \lambda_{tgi j}^{\text{EN}} E_{gij}$$

$$\log(\lambda_{tgi j}^{\text{AR}}) = \mu^{\text{AR}} + \alpha_j^{\text{AR}} + \beta_g^{\text{AR}}$$

$$\log(\lambda_{tgi j}^{\text{NE}}) = \mu^{\text{NE}} + \alpha_j^{\text{NE}} + \beta_g^{\text{NE}}$$

$$\begin{aligned} \log(\lambda_{tgi j}^{\text{EN}}) = & \mu^{\text{EN}} + \alpha_j^{\text{EN}} + \beta_g^{\text{EN}} + \gamma^{\text{EN}} t + \delta^{\text{EN}} \sin\left(\frac{t}{52} 2\pi\right) \\ & + \varepsilon^{\text{EN}} \cos\left(\frac{t}{52} 2\pi\right) \end{aligned}$$

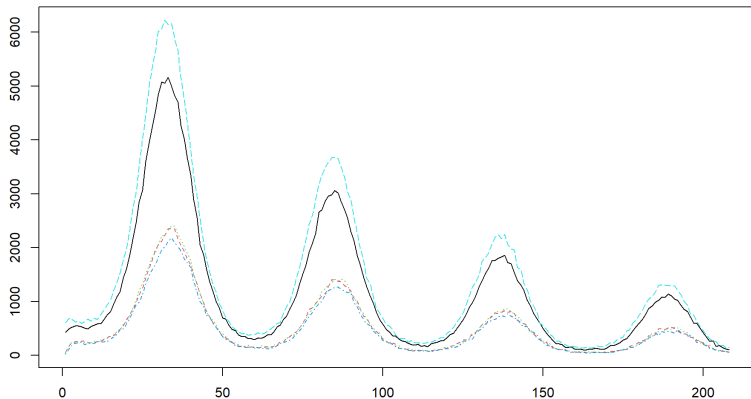
$$\text{logit}(p_{tgi j}) = \mu^{(\text{p})} + \alpha_j^{(\text{p})} + \beta_g^{(\text{p})}$$

# Simulation Study Overview (Work-In-Progress)

- Data generated according to the Latent Model in R
- Model fitting performed via HMC with Stan
- Scenarios to Consider:
  - Similar EE-type model, ignoring missingness entirely
  - Similar EE-type model, imputation via statistical model
  - Similar EE-type model, imputation via ML or MICE
  - TSIR models (i.e. Neg. Bin. likelihood) for each of the above

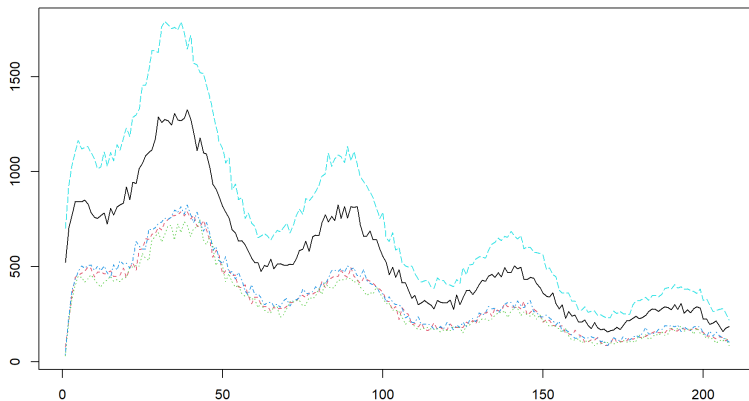
# Simulated Data (Poisson)

# Cases over time, by race

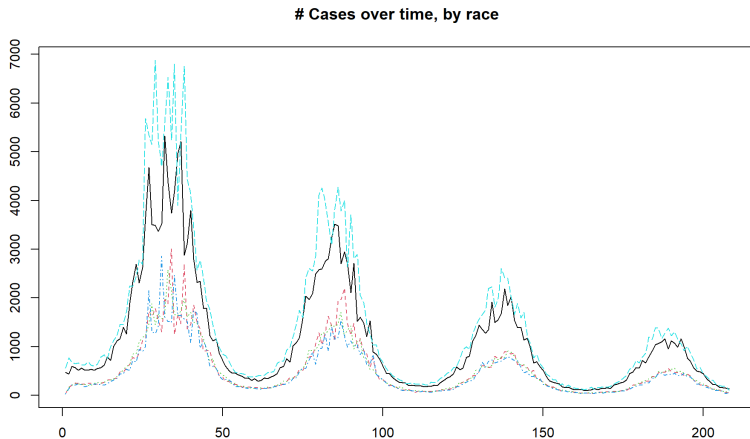


# Simulated Data (Poisson)

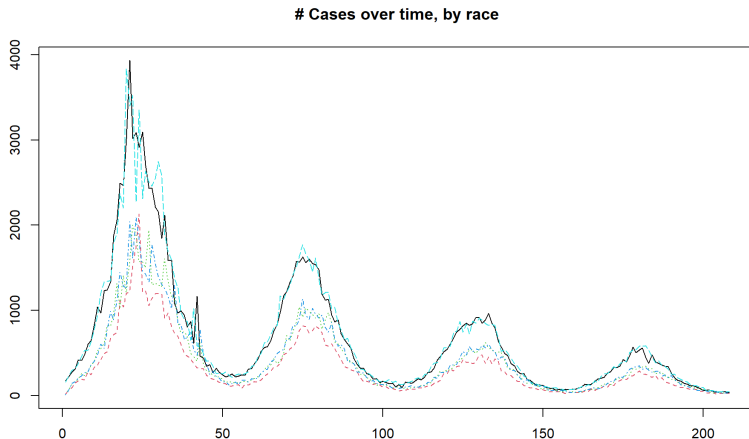
# Cases over time, by race



# Simulated Data (Negative Binomial)



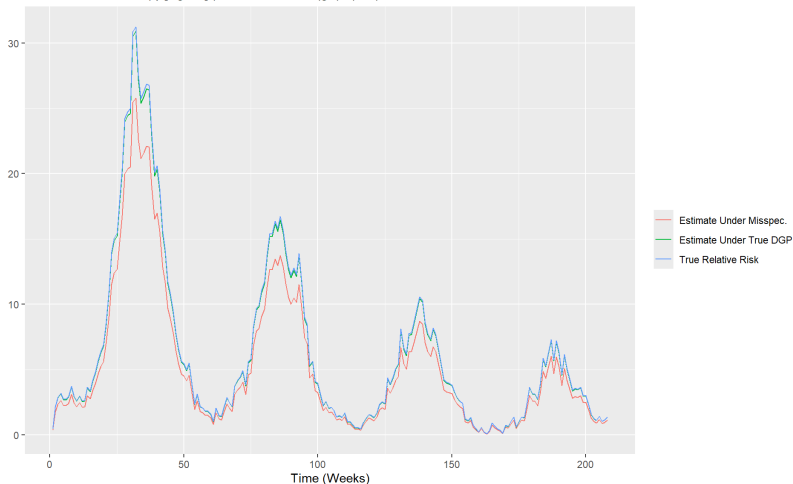
# Simulated Data (Negative Binomial)



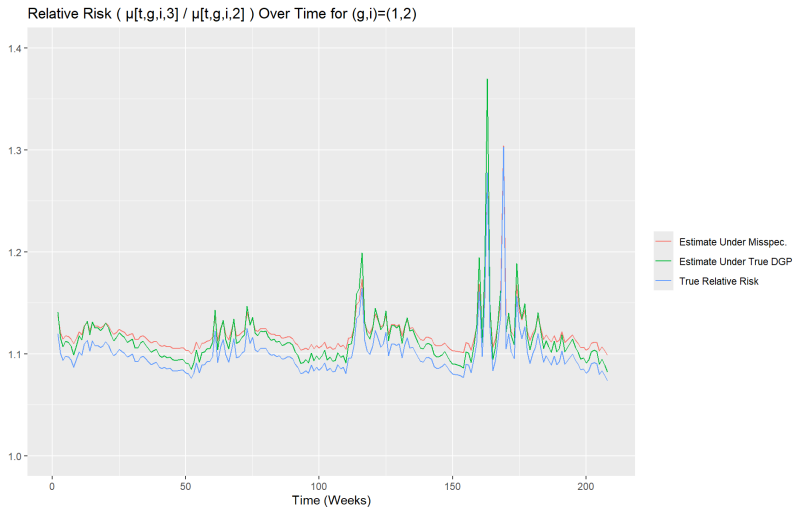


# Preliminary Results (Idealized Scenario)

Total Hazard Rate (  $\mu[t,g,i,3]$  ) Over Time for  $(g,i)=(1,2)$

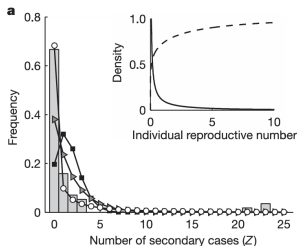


# Preliminary Results (Idealized Scenario)



# Motivation for Zelner Contact-Heterogeneity Model

- Observation: Individual-level contact-tracing data seems to suggest the average number of secondary-infections caused by an individual (i.e.  $R_0$ ) exhibits individual-level heterogeneity (see Lloyd-Smith et al. 2005)



- Observation: If we wish to stratify a disease model by geography and demography, there is no obvious way in the aforementioned frameworks to specify how new infections are doled-out *across* strata

# Derivation #4: Zelner Contact-Heterogeneity Model

- (1) Suppose the (latent) individual  $R_0$  of infected  $i$  in geography  $g$  during the interval  $(t-1, t)$  is given by

$$r_{tgi} \stackrel{\text{iid}}{\sim} \text{Gamma}\left(\underbrace{\frac{R_0}{\theta}}_{\text{shape}}, \underbrace{\theta}_{\text{scale}}\right)$$

- (2) The total (latent) infectiousness at time  $t$  for geography  $g$  becomes

$$(r_{t\bullet} \mid Y_{tg} = y_{tg}) = \left(\sum_{i=1}^{y_{tg}} r_{ti}\right) \sim \text{Gamma}\left(\frac{R_0}{\theta} y_{tg}, \theta\right)$$

- (3) Assume homogeneous mixing, and that all of the latent infectiousness is deposited in the single time period after infection. Then the force of infection is

$$\lambda_{tg} = \zeta \frac{r_{(t-1)g}}{n_g} + (1 - \zeta) \sum_{g' \neq g} \frac{r_{(t-1)g'}}{N - n_{g'}}$$

- (4) Following a process almost identical to the EE model framework:

$$(Y_{tg} \mid \mathbf{r}_{(t-1)}) \sim \text{Pois}(n_g \lambda_{tg})$$

# Recap

- Going forward, achieving health equity will require more sophisticated methods for handling MNAR data

# Recap

- Going forward, achieving health equity will require more sophisticated methods for handling MNAR data
- Our attempt at this is to augment the EE-like and TSIR-like models of disease incidence with a selection-model component

# Recap

- Going forward, achieving health equity will require more sophisticated methods for handling MNAR data
- Our attempt at this is to augment the EE-like and TSIR-like models of disease incidence with a selection-model component
- Demonstrating efficacy will revolve around comparing relative risk estimates between different models and missing-data techniques

# Recap

- Going forward, achieving health equity will require more sophisticated methods for handling MNAR data
- Our attempt at this is to augment the EE-like and TSIR-like models of disease incidence with a selection-model component
- Demonstrating efficacy will revolve around comparing relative risk estimates between different models and missing-data techniques
- Future work includes



# Recap

- Going forward, achieving health equity will require more sophisticated methods for handling MNAR data
- Our attempt at this is to augment the EE-like and TSIR-like models of disease incidence with a selection-model component
- Demonstrating efficacy will revolve around comparing relative risk estimates between different models and missing-data techniques
- Future work includes
  - Determining parameterizations that dole-out infectiousness among strata in meaningful ways

# Recap

- Going forward, achieving health equity will require more sophisticated methods for handling MNAR data
- Our attempt at this is to augment the EE-like and TSIR-like models of disease incidence with a selection-model component
- Demonstrating efficacy will revolve around comparing relative risk estimates between different models and missing-data techniques
- Future work includes
  - Determining parameterizations that dole-out infectiousness among strata in meaningful ways
  - Determining conditions for local and global identifiability

# Recap

- Going forward, achieving health equity will require more sophisticated methods for handling MNAR data
- Our attempt at this is to augment the EE-like and TSIR-like models of disease incidence with a selection-model component
- Demonstrating efficacy will revolve around comparing relative risk estimates between different models and missing-data techniques
- Future work includes
  - Determining parameterizations that dole-out infectiousness among strata in meaningful ways
  - Determining conditions for local and global identifiability
  - Determining validity of rare-disease assumption and consequences of violation

# References I



Bartlett, M. S. (1956). "Deterministic and Stochastic Models for Recurrent Epidemics". In: *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability, Volume 4: Contributions to Biology and Problems of Health*. Vol. 3.4. University of California Press, pp. 81–110. (Visited on 11/03/2024).



Bauer, Cici and Jon Wakefield (2018). "Stratified space–time infectious disease modelling, with an application to hand, foot and mouth disease in China". In: *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 67.5, pp. 1379–1398. ISSN: 1467-9876. DOI: 10.1111/rssc.12284. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1111/rssc.12284>.



Bjørnstad, Ottar N., Bärbel F. Finkenstädt, and Bryan T. Grenfell (2002). "Dynamics of Measles Epidemics: Estimating Scaling of Transmission Rates Using a Time Series SIR Model". In: *Ecological Monographs* 72.2, pp. 169–184. ISSN: 1557-7015. DOI: 10.1890/0012-9615(2002)072[0169:DOMEES]2.0.CO;2.

# References II



Held, Leonhard, Michael Höhle, and Mathias Hofmann (2005). "A statistical framework for the analysis of multivariate infectious disease surveillance counts". In: *Statistical Modelling* 5.3, pp. 187–199. ISSN: 1471-082X. DOI: 10.1191/1471082X05st098oa. URL: <https://doi.org/10.1191/1471082X05st098oa>.



Held, Leonhard and Michaela Paul (2012). "Modeling seasonality in space-time infectious disease surveillance data". In: *Biometrical Journal* 54.6, pp. 824–843. ISSN: 1521-4036. DOI: 10.1002/bimj.201200037. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/bimj.201200037>.



Kendall, David G. (1949). "Stochastic Processes and Population Growth". In: *Journal of the Royal Statistical Society: Series B (Methodological)* 11.2, pp. 230–264. ISSN: 0035-9246. DOI: 10.1111/j.2517-6161.1949.tb00032.x. (Visited on 11/03/2024).

# References III



Lloyd-Smith, J. O. et al. (2005). “Superspreading and the effect of individual variation on disease emergence”. In: *Nature* 438.7066, pp. 355–359. ISSN: 1476-4687. DOI: 10.1038/nature04153. URL: <https://doi.org/10.1038/nature04153>.



Trangucci, Rob, Yang Chen, and Jon Zelner (2023). “Modeling racial/ethnic differences in COVID-19 incidence with covariates subject to nonrandom missingness”. In: *The Annals of Applied Statistics* 17.4, pp. 2723–2758. ISSN: 1932-6157, 1941-7330. DOI: 10.1214/22-AOAS1711. URL: <https://projecteuclid.org/journals/annals-of-applied-statistics/volume-17/issue-4/Modeling-racial-ethnic-differences-in-COVID-19-incidence-with-covariates/10.1214/22-AOAS1711.full> (visited on 11/04/2024).



Wakefield, Jon, Tracy Qi Dong, and Vladimir N. Minin (2019). “Spatio-Temporal Analysis of Surveillance Data”. In: *Handbook of Infectious Disease Data Analysis*. 1st. Chapman and Hall/CRC, pp. 455–475. ISBN: 978-1-315-22291-2.

# References IV



Zelner, Jon et al. (2020). "Understanding the Importance of Contact Heterogeneity and Variable Infectiousness in the Dynamics of a Large Norovirus Outbreak". In: *Clinical Infectious Diseases* 70.3, pp. 493–500. ISSN: 1058-4838. DOI: [10.1093/cid/ciz220](https://doi.org/10.1093/cid/ciz220).