

SIS models for recurrent infections

SISMID/July 11–13, 2018

Instructors: Vladimir Minin, Kari Auranen, Elizabeth Halloran

Outline

- ▶ Recurrent infections
- ▶ A simple SIS model without transmission
 - ▶ Complete-data likelihood
- ▶ Model extensions
 - ▶ Several subtypes of a pathogen
 - ▶ Transmission in small mixing groups
- ▶ Incomplete observations
 - ▶ Discrete-time Markov transition models
 - ▶ Continuous-time Markov processes with Bayesian data augmentation and reversible jump MCMC
- ▶ A computer class exercise of an SIS model without transmission and with completely observed data

Background

- ▶ Many infections can be considered recurrent, i.e., occurring as an alternating series of presence and absence of infection
 - ▶ Nasopharyngeal carriage of *Streptococcus pneumoniae*
(Auranen et al.; Cauchemez et al.; Melegaro et al.)
 - ▶ Nasopharyngeal carriage of *Neisseria meningitidis*
 - ▶ multi-resistant *Staphylococcus aureus* (Cooper et al.)
 - ▶ HPV (human papilloma virus) infection
 - ▶ some parasitic infections (e.g. Nagelkerke et al.)
- ▶ Many of the above infections are asymptomatic, which means that observation requires active sampling of the underlying epidemiological states
- ▶ Exact acquisition and clearance times often remain unobserved \Rightarrow incompletely observed data

A binary Markov process

A simple model for a recurrent infection is the binary Markov process:

- ▶ The state of the individual alternates between “susceptible” (state 0) and “infected” (state 1)
- ▶ Hazard of acquiring infection is β :

$$P(\text{acquisition in } [t, t + dt[\mid \text{susceptible at time } t-) \simeq \beta dt$$

- ▶ Hazard of clearing infection is μ :

$$P(\text{clearance in } [t, t + dt[\mid \text{infected at time } t-) \simeq \mu dt$$

Complete data

- ▶ For each individual i , the complete data include the times of acquisition and clearance during the observation period $[0, T]$:
 - ▶ Denote the ordered acquisition times for individual i during $]0, T[$ by $\mathbf{t}^{(i)} = (t_{i1}, \dots, t_{iN_{01}^{(i)}})$
 - ▶ Denote the ordered clearance times for individual i during $]0, T[$ by $\mathbf{r}^{(i)} = (r_{i1}, \dots, r_{iN_{10}^{(i)}})$
 - ▶ Denote the ordered acquisition and clearance times together as $u_{i1} = 0, u_{i2}, \dots, u_{i,N^{(i)}} = T$
 - ▶ Note: these include times 0 and T
(so that $N^{(i)} = N_{01}^{(i)} + N_{10}^{(i)} + 2$)

Keeping track who is susceptible

- ▶ The binary indicators for individual i to be susceptible or infected at time t are denoted by $Y_0^{(i)}(t)$ and $Y_1^{(i)}(t)$, respectively
 - ▶ For the simple binary model, $Y_1^{(i)}(t) = 1 - Y_0^{(i)}(t)$ for all times $t \geq 0$, i.e. the individual is always either susceptible or infected
 - ▶ Both indicators are taken to be *predictable*, i.e., their values at time t are determined by their initial values and the complete data observed up to time t —
 - ▶ In practice, this means that the values of $Y_0^{(i)}(t)$ and $Y_1^{(i)}(t)$ can be calculated from the observed data and these indicators can be used as shorthand when writing the likelihood function

Process of acquisitions

- ▶ In each individual, acquisitions (i.e. infections) occur with intensity $\beta Y_0^{(i)}(t)$
 - ▶ The intensity is β when the individual is in state 0 (susceptible) and 0 when the individual is in state 1 (infected)
- ▶ The probability density of the acquisition events is proportional to

$$\prod_{k=1}^{N^{(i)}} \left[\beta \mathbf{1}(u_k \text{ is time of acq.}) e^{-\beta Y_0^{(i)}(u_k)(u_k - u_{k-1})} \right]$$

total time susceptible

$$\propto \beta^{N_{01}^{(i)}} \times \exp \left\{ -\beta \overbrace{\sum_{k=1}^{N^{(i)}} Y_0^{(i)}(u_k)(u_k - u_{k-1})}^{\text{total time susceptible}} \right\}$$

Process of clearances

- ▶ In each individual, clearances occur with intensity $\mu Y_1^{(i)}(t)$
 - ▶ The intensity is μ when the individual is in state 1 (infected) and 0 when the individual is in state 0 (susceptible)
- ▶ The probability density of the clearance events is proportional to

$$\prod_{k=1}^{N^{(i)}} \left[\mu^{1(u_k \text{ is time of clearance})} e^{-\mu Y_1^{(i)}(u_k)(u_k - u_{k-1})} \right]$$
$$= \mu^{N_{10}^{(i)}} \times \exp \left\{ -\mu \overbrace{\sum_{k=1}^{N^{(i)}} Y_1^{(i)}(u_k)(u_k - u_{k-1})}^{\text{total time infected}} \right\}$$

Complete data likelihood

- The likelihood function of parameters β and μ , based on the complete data from individual i :

$$\begin{aligned} & \overbrace{f(\mathbf{t}^{(i)}, \mathbf{r}^{(i)} | \beta, \mu)} \\ & L_i(\beta, \mu; \mathbf{t}^{(i)}, \mathbf{r}^{(i)}) \\ &= \beta^{N_{01}^{(i)}} \mu^{N_{10}^{(i)}} \times e^{-\sum_{k=1}^{N^{(i)}} (\beta Y_0^{(i)}(u_k) + \mu Y_1^{(i)}(u_k))(u_k - u_{k-1})} \\ &= \beta^{N_{01}^{(i)}} \mu^{N_{10}^{(i)}} \times \exp\left(-\int_0^T \{\beta Y_0^{(i)}(u) + \mu Y_1^{(i)}(u)\} du\right) \end{aligned}$$

- Likelihood based on all M individuals is $\prod_{i=1}^M L_i(\beta, \mu; \mathbf{t}^{(i)}, \mathbf{r}^{(i)})$

Model extension 1: Several subtypes of infection

- ▶ If the pathogen has K subtypes/strains, infections with each can be modelled as separate states
- ▶ Let the hazards from making a transition from state s to state r be α_{sr} , $s, r \in \{0, \dots, K\}$
 - ▶ This means that infection for a susceptible occurs with rate α_{0r} , $r = 1, \dots, K$
 - ▶ Also direct transitions from infection with type s to infection with type r are possible if we allow $\alpha_{sr} > 0$ also when $s > 0$
- ▶ Transitions for each individual are now modelled as follows:

$$P(\text{individual } i \text{ makes transition from } s \text{ to } r \text{ in } [t, t + dt) \simeq \alpha_{sr} Y_s^{(i)}(t) dt$$

Model extension 2: Modelling transmission

- ▶ The hazard of infection may depend on the presence of infected individuals in the family, day care group, school class etc.
 - ▶ The statistical unit is defined by the relevant mixing group
- ▶ Denote $H_t^{(i, fam)}$ the joint infection status of all members in the mixing group (e.g. family) of individual i at time t —
- ▶ For a single-type pathogen, infections are modeled as follows:

$$P(\text{infection for } i \text{ in } [t, t + dt] | H_{t-}^{(i, fam)}) \simeq \alpha_{01}^{(i)}(t) Y_0^{(i)}(t) dt \equiv \frac{\beta C^{(i)}(t)}{M_{fam}^{(i)} - 1} Y_0^{(i)}(t) dt$$

where $C^{(i)}(t)$ is the number of infected individuals in i 's family (of size $M_{fam}^{(i)}$) at time t —; note that $C^{(i)}(t)$ can be calculated from the state-indicator variables of i 's family members

Complete data likelihood: the general expression

- ▶ For M individuals followed from time 0 to time T , the *complete data* record all transitions between states s and r :

$$\mathbf{x}_{\text{complete}} = \{T_{sr}^{(ik)}; s, r = 0, 1 (s \neq r), k = 1, \dots, N_{sr}^{(i)}(T), i = 1, \dots, M\}$$

- ▶ The likelihood of the rate parameters θ , based on the complete (event-history) data

$$\underbrace{f(\mathbf{x}_{\text{complete}}|\theta)}_{L(\theta; \mathbf{x}_{\text{complete}})} = \prod_i^N \prod_{r \neq s} \prod_k^{N_{sr}^{(i)}(T)} \left[\alpha_{sr}^{(i)}(T_{sr}^{(ik)}) \times \exp \left(- \int_0^T \alpha_{sr}^{(i)}(u) Y_s^{(i)}(u) du \right) \right]$$

Remarks

- ▶ Although the likelihood expression was written as a product of individual likelihood contributions, it is valid even if the individual processes are dependent on the infection outcomes of *other* individuals (as when modeling transmission)
- ▶ The likelihood is correctly normalized with respect to any number of events occurring between times 0 and T
 - ▶ This is crucial when performing MCMC computations through data augmentation with an unknown number of events
- ▶ These results are somewhat non-trivial and require the theory of counting processes (Andersen et al.)

Incomplete observations

- ▶ Usually we do not observe complete data (= all infection and clearance times)
- ▶ Instead, the status (infection stage) $X_j^{(i)}$ of each individual is observed at pre-defined times $t_j^{(i)}$
 - ▶ This creates *incomplete data*: the process is only observed at discrete times (panel data)
 - ▶ The observed data likelihood is now a complicated function of the model parameters
- ▶ How to estimate the underlying continuous process from discrete observations?
 - ▶ (A) Discrete-time Markov transition model
 - ▶ (B) Continuous-time Markov transition model with Bayesian data augmentation

(A) Markov transition models

- ▶ Treat the problem as a discrete-time Markov transition model
- ▶ For simplicity, assume equal-length ($= \Delta$) time steps indexed by $t = 1, 2, 3, \dots$, and time-homogeneous one-time-step transition probabilities:

$$p_{sr} \equiv P(X_{t+1}^{(i)} = r | X_t^{(i)} = s) \quad \text{for } t = 1, 2, \dots; s, r = 0, \dots, K$$

- ▶ This defines matrix P_Δ of one-step transition probabilities with entries $[P_\Delta]_{sr} = p_{sr}$, where $\sum_r p_{sr} = 1$ for each $s = 0, \dots, K$

Likelihood function under the discrete model

- ▶ Denote the observed numbers of one-step transitions in all study subjects by N_{sr} , $s, r = 0, \dots, K$
- ▶ The likelihood of the unknown transition probabilities (i.e. the elements of matrix P_Δ) is now particularly simple:

$$L(P_\Delta) = \prod_{s,r} [p_{sr}(\Delta)]^{N_{sr}(T)} = \prod_{s,r} [P_\Delta]_{sr}^{N_{sr}(T)}$$

- ▶ When observation are actually made at intervals $k\Delta$ (e.g. $\Delta =$ day and $k = 28$), the likelihood is

$$L(P_\Delta) = \prod_{s,r} [P_\Delta^k]_{sr}^{N_{sr}(T)}$$

Modeling transmission

- ▶ In a mixing group of size M , the state space is $\chi \times \chi \times \dots \chi$, where χ is the state space for one individual
 - ▶ In a binary infection model, the individual state space is $\{0, 1\}$
 - ▶ In a family of three individuals the state space is then $\{(0, 0, 0), (1, 0, 0), (0, 1, 0), (0, 0, 1), (1, 1, 0), (1, 0, 1), (0, 1, 1), (1, 1, 1)\}$
 - ▶ For M individuals, the dimension of the state space is 2^M
- ▶ Application to pneumococcal carriage in families (Melegaro et al.)
 - ▶ The transition probability matrix in a family of 3 (next page), assuming the same probabilities (per day) for each family member
 - ▶ Notation: $q_{ii} = 1 - (\text{sum of the non-diagonal elements on the } i\text{th row})$

Transition probability matrix

$$P_{\Delta} = \begin{matrix} & \begin{matrix} (0,0,0) & (1,0,0) & (0,1,0) & (0,0,1) & (1,1,0) & (1,0,1) & (0,1,1) & (1,1,1) \end{matrix} \\ \begin{pmatrix} q_{11} & \kappa & \kappa & \kappa & 0 & 0 & 0 & 0 \\ \mu & q_{22} & 0 & 0 & \beta/2 + \kappa & \beta/2 + \kappa & 0 & 0 \\ \mu & 0 & q_{33} & 0 & \beta/2 + \kappa & 0 & \beta/2 + \kappa & 0 \\ \mu & 0 & 0 & q_{44} & 0 & \beta/2 + \kappa & \beta/2 + \kappa & 0 \\ 0 & \mu & \mu & 0 & q_{55} & 0 & 0 & \beta + \kappa \\ 0 & \mu & 0 & \mu & 0 & q_{66} & 0 & \beta + \kappa \\ 0 & 0 & \mu & \mu & 0 & 0 & q_{77} & \beta + \kappa \\ 0 & 0 & 0 & 0 & \mu & \mu & \mu & q_{88} \end{pmatrix} \end{matrix}$$

Potential problems

- ▶ The dimension of the state space
 - ▶ With M individuals and $K + 1$ types of infection, the dimension of the state space is $(K + 1)^M$
 - ▶ With 13 serotypes and 25 individuals (see Hoti et al.), the dimension is $\sim 4.5 \times 10^{28}$
- ▶ Non-Markovian sojourn times
 - ▶ e.g. a Weibull duration of infection may be more realistic than the exponential one
- ▶ Handling of varying observation intervals and individuals with completely missing data are still cumbersome

(B) Bayesian data augmentation

- ▶ If we retain the continuous-time model formulation, unobserved event times of acquisition and clearance can be taken as additional model unknowns (parameters)
- ▶ Statistical inference on all model unknowns (parameters θ and event times x_{complete})

$$\overbrace{f(x_{\text{observed}} | x_{\text{complete}})}^{\text{observation model}} \quad \overbrace{f(x_{\text{complete}} | \theta)}^{\text{complete data likelihood}} \quad \overbrace{f(\theta)}^{\text{prior}}$$

- ▶ The observed data x_{observed} contain only the current status of infection in each study subject at predefined time points
- ▶ The observation model often only ensures agreement with the observed data (as an indicator function)
- ▶ The computational problem:
how to sample from $f(x_{\text{complete}} | x_{\text{observed}}, \theta)$?

Sampling algorithm

- ▶ Initialize the model parameters and the latent processes (i.e. the unobserved event times)
- ▶ For each individual, update the unobserved event times
 - ▶ Update the current iterates of the event times using standard MH
 - ▶ Add/delete episodes of infection and non-infection using reversible jump MH
 - ▶ with 0.5 probability propose to add a new episode
 - ▶ with 0.5 probability propose to delete an existing episode
- ▶ Update the model parameters using single-step MH
- ▶ Iterate the above updating steps for a given number of MCMC iterations

Adding/deleting episodes

- ▶ Choose one interval at random from among the K sampling intervals (see page+2)
- ▶ Choose to add an episode (delete an existing episode) within the chosen interval with probability $\pi_{\text{add}} = 0.5$ ($\pi_{\text{delete}} = 0.5$)
 - ▶ If 'add', choose random event times $\bar{t}_1 < \bar{t}_2$ uniformly from Δ (= the length of the sampling interval). These define the new episode.
 - ▶ If 'delete', delete the two event times
- ▶ The 'add' move is accepted with probability (Metropolis-Hastings acceptance ratio)

$$\min \left(\frac{f(x_{\text{observed}} | x_{\text{complete}}^*) f(x_{\text{complete}}^* | \theta) q(x_{\text{complete}} | x_{\text{complete}}^*)}{f(x_{\text{observed}} | x_{\text{complete}}) f(x_{\text{complete}} | \theta) q(x_{\text{complete}}^* | x_{\text{complete}})}, 1 \right)$$

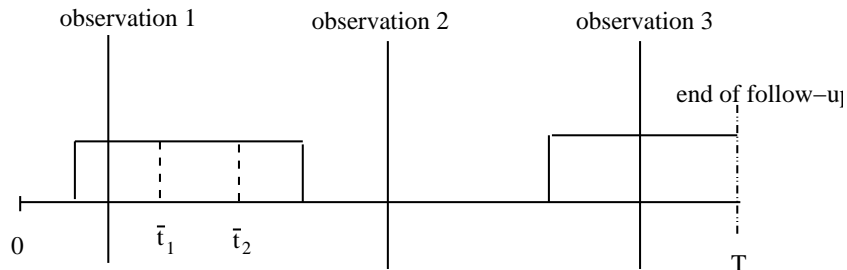
Adding/deleting episodes cont.

- ▶ The ratio of the proposal densities is

$$\frac{q(x_{\text{complete}} | x_{\text{complete}}^*)}{q(x_{\text{complete}}^* | x_{\text{complete}})} = \frac{\pi_{\text{delete}} \frac{1}{K} \frac{1}{L}}{\pi_{\text{add}} \frac{1}{K} \frac{1}{L} \frac{2}{\Delta^2}} = \frac{\Delta^2}{2}$$

- ▶ The ratio of the proposal densities in the 'delete' move is the inverse of the expression above
- ▶ Technically, the add/delete step relies on so called reversible jump MCMC (see page+2)
- ▶ Reversible jump types should be devised to assure irreducibility of the Markov chain
- ▶ For a more complex example, see e.g. Hoti et al.

Adding/deleting latent processes cont.



The number of sampling intervals $K=4$

The number of 'sub-episodes' within the second interval $L=2$

Reversible jump MCMC

- ▶ “When the number of things you don’t know is one of the things you don’t know”
- ▶ For example, under incomplete observation of the previous (Markov) processes, the exact number of events is not observed
- ▶ This requires a joint model over ‘sub-spaces’ of different dimensions
- ▶ And a method to do numerical integration (MCMC sampling) in the joint state space

References

- [1] Andersen et al. "Statistical models based on counting processes", Springer, 1993
- [2] Auranen et al. "Transmission of pneumococcal carriage in families – a latent Markov process model for binary data. J Am Stat Assoc 2000; 95:1044-1053.
- [3] Melegaro et al. Estimating the transmission parameters of pneumococcal carriage in families. Epidemiol Infect 2004; 132:433-441.
- [4] Cauchemez et al. Streptococcus pneumoniae transmission according to inclusion in conjugate vaccines: Bayesian analysis of a longitudinal follow-up in schools. BMC Infectious Diseases 2006, 6:14.
- [5] Nakelkerke et al. Estimation of parasitic infection dynamics when detectability is imperfect. Stat Med 1990; 9:1211-1219.
- [6] Cooper et al. "An augmented data method for the analysis of nosocomial infection data. Am J Epidemiol 2004; 168:548-557.
- [7] Bladt et al. "Statistical inference for discretely observed Markov jump processes. J R Statist Soc B 2005; 67:395-410.
- [8] Andersen et al. Multi-state models for event history analysis. Stat Meth Med Res 2002; 11:91-115.
- [9] Hoti et al. Outbreaks of Streptococcus pneumoniae carriage in day care cohorts in Finland – implications to elimination of carriage. BMC Infectious Diseases, 2009 (in press)
- [10] Green P. Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. Biometrika 1995; 82:711-732.