# BAYESIAN MODEL SELECTION FOR PARTIALLY OBSERVED EPIDEMIC MODELS

**Panayiota Touloupou**

joint work with Simon E.F. Spencer, Bärbel Finkenstädt
Rand, Peter Neal, Trevelyan J. McKinley

Warwick
**Statistics**

CRiSM Workshop: Estimating Constants
April 21, 2016

# OUTLINE

1 MOTIVATION

2 METHODS

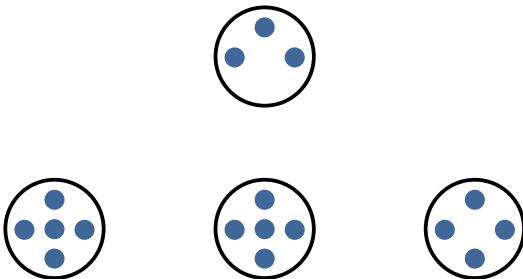3 SIMULATION STUDIES

4 REAL DATA ANALYSIS

5 CONCLUSIONS

**MOTIVATION**
000

**METHODS**
00000000000

**SIMULATION STUDIES**
000000

**REAL DATA ANALYSIS**
0000000000

**CONCLUSIONS**
00

## OUTLINE

1 **MOTIVATION**

2 **METHODS**

3 **SIMULATION STUDIES**

4 **REAL DATA ANALYSIS**

5 **CONCLUSIONS**

MOTIVATION
●○○

METHODS
○○○○○○○○○○○

SIMULATION STUDIES
○○○○○○

REAL DATA ANALYSIS
○○○○○○○○○○

CONCLUSIONS
○○

# STATISTICAL EPIDEMIC MODELLING

- Insights into dynamics of infectious diseases
  - ➣ Prevention
  - ➣ Control spread of the disease

- Epidemiological data present several challenges
  - ➣ Missing data (typically high dimensional)
  - ➣ Diagnostic tests imperfect

- Model selection
  - ➣ Each model an epidemiologically important hypothesis
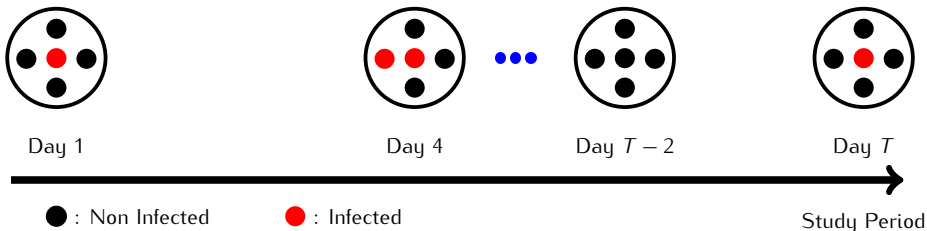
# OUR SETUP

- Longitudinal observations
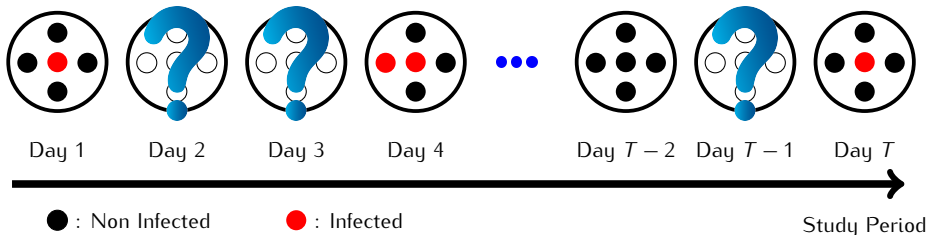- Individuals form groups (e.g. households)



● : Individual

# OUR SETUP

- Longitudinal observations
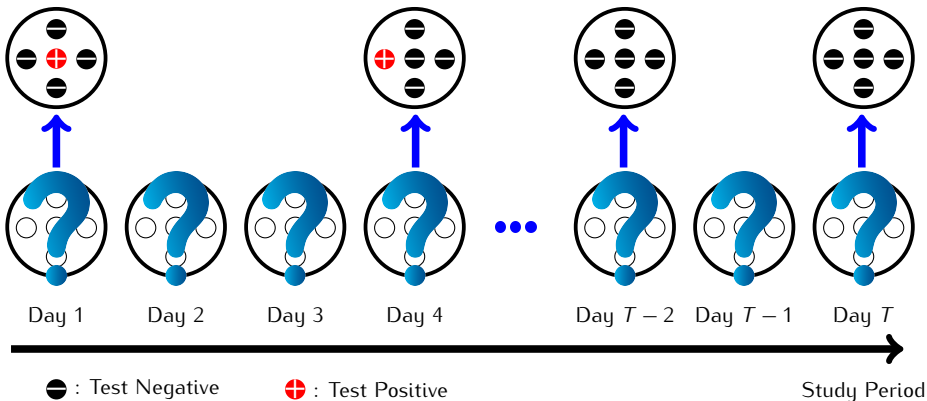- Individuals form groups (e.g. households)



Day 1          Day 4          Day $T-2$          Day $T$

● : Non Infected          ● : Infected

Study Period

# OUR SETUP

- Longitudinal observations
- Individuals form groups (e.g. households)



- Day 1   Day 2   Day 3   Day 4   •••   Day $T-2$   Day $T-1$   Day $T$

● : Non Infected       ● : Infected

Study Period

# OUR SETUP

- Longitudinal observations
- Individuals form groups (e.g. households)



⊖ : Test Negative    ⊕ : Test Positive    Study Period

## OBJECTIVE

- Analysis of this type of data can be challenging
  - ➤ Times of acquiring and clearing infection are unobserved
  - ➤ Intractable likelihood – need to know missing times
  - ➤ Usual solution: large scale data augmentation MCMC

- **Bayesian model selection**
  - ➤ Evidence in favour of candidate models
  - ➤ Each model an epidemiologically important hypothesis

<table>
<tr>
<td rowspan="2">OBJECTIVES:</td>
<td>● Develop statistical tools for comparison of competing hypotheses</td>
</tr>
<tr>
<td>● Special attention on missing data</td>
</tr>
</table>

# OUTLINE

1. **MOTIVATION**

2. **METHODS**

3. **SIMULATION STUDIES**

4. **REAL DATA ANALYSIS**

5. **CONCLUSIONS**

MOTIVATION
000

METHODS
●○○○○○○○○○○

SIMULATION STUDIES
000000

REAL DATA ANALYSIS
0000000000

CONCLUSIONS
00

# MODEL SELECTION FOR EPIDEMICS

**A lot of epidemiologically interesting questions take the form of model selection questions**

- What is the transmission mechanism of the disease?

- Do individuals develop immunity over time?

- Do water troughs spread *E. coli* O157?

# POSTERIOR PROBABILITIES AND MARGINAL LIKELIHOODS

- Would like the posterior probability in favour of model $i$

$$P(M_i|\boldsymbol{y}) = \frac{\pi(\boldsymbol{y}|M_i)P(M_i)}{\sum_j \pi(\boldsymbol{y}|M_j)P(M_j)}$$

# POSTERIOR PROBABILITIES AND MARGINAL LIKELIHOODS

- Would like the posterior probability in favour of model $i$

$$P(M_i|\boldsymbol{y}) = \frac{\pi(\boldsymbol{y}|M_i)P(M_i)}{\sum_j \pi(\boldsymbol{y}|M_j)P(M_j)}$$

- Equivalently, the Bayes factor comparing models $i$ and $j$

$$B_{ij} = \frac{\pi(\boldsymbol{y}|M_i)}{\pi(\boldsymbol{y}|M_j)}$$

# POSTERIOR PROBABILITIES AND MARGINAL LIKELIHOODS

- Would like the posterior probability in favour of model $i$

$$P(M_i|\boldsymbol{y}) = \frac{\pi(\boldsymbol{y}|M_i)P(M_i)}{\sum_j \pi(\boldsymbol{y}|M_j)P(M_j)}$$

- Equivalently, the Bayes factor comparing models $i$ and $j$

$$B_{ij} = \frac{\pi(\boldsymbol{y}|M_i)}{\pi(\boldsymbol{y}|M_j)}$$

- All we need is the marginal likelihood,

$$\pi(\boldsymbol{y}|M_i) = \int \pi(\boldsymbol{y}|\boldsymbol{\theta}, M_i)\pi(\boldsymbol{\theta}|M_i)\,\mathrm{d}\boldsymbol{\theta}$$

but how can we calculate it?

MOTIVATION
000

METHODS
0000000000

SIMULATION STUDIES
000000

REAL DATA ANALYSIS
0000000000

CONCLUSIONS
00

# MARGINAL LIKELIHOOD ESTIMATION

- Most direct approach: **Importance sampling**
  - ➤ Use asymptotic normality of the posterior to find efficient proposal

- Many existing other approaches:
  - ➤ Harmonic mean
  - ➤ Chib's methods
  - ➤ Power posteriors
  - ➤ Bridge sampling

MOTIVATION
000

METHODS
0000●000000

SIMULATION STUDIES
000000

REAL DATA ANALYSIS
0000000000

CONCLUSIONS
00

# IMPORTANCE SAMPLING[1]

1. Obtain samples from the posterior $\pi(\boldsymbol{\theta}|\boldsymbol{y})$ with MCMC

2. Use MCMC samples to inform the proposal distribution $\Rightarrow q(\boldsymbol{\theta})$

3. Draw $N$ samples from $q(\boldsymbol{\theta})$

4. Estimate the marginal likelihood by

$$\widehat{\pi}_{IS}(\boldsymbol{y}) = \sum_{i=1}^{N} \frac{\pi(\boldsymbol{y}|\boldsymbol{\theta}_i)\pi(\boldsymbol{\theta}_i)}{q(\boldsymbol{\theta}_i)}$$

---

[1]Clyde et al. (2007). Current Challenges in Bayesian Model Choice

# MISSING DATA!

But how to deal with the missing data?

# MISSING DATA!

But how to deal with the missing data?

MOTIVATION
○○○

METHODS
○○○○○●○○○○○

SIMULATION STUDIES
○○○○○○

REAL DATA ANALYSIS
○○○○○○○○○○

CONCLUSIONS
○○

# IMPORTANCE SAMPLING WITH MISSING DATA

1. Obtain samples from the joint posterior $\pi(x, \theta|y)$ with MCMC

2. Use MCMC samples to inform the proposal distribution $\Rightarrow q(\theta)$

3. Draw $N$ samples from $q(\theta)$

4. For each sampled $\theta_i$ draw missing data $x_i$ from the full conditional using Forward Filtering Backward Sampling

5. Estimate the marginal likelihood by

$$\widehat{\pi}_{IS}(y) = \sum_{i=1}^{N} \frac{\pi(y|x_i, \theta_i)\ \pi(x_i|\theta_i)\ \pi(\theta_i)}{\pi(x_i|y, \theta_i)\ q(\theta_i)}$$

MOTIVATION
000

METHODS
0000000●0000

SIMULATION STUDIES
000000

REAL DATA ANALYSIS
0000000000

CONCLUSIONS
00

# HARMONIC MEAN[2]

- The marginal likelihood $\pi(\boldsymbol{y})$ can be approximated

$$\widehat{\pi}_{HM}(\boldsymbol{y}) = \left[ \frac{1}{N} \sum_{i=1}^{N} \frac{1}{\pi(\boldsymbol{y}|\boldsymbol{x}_i, \boldsymbol{\theta}_i)} \right]^{-1}$$

  based on $N$ draws $(\boldsymbol{x}_1, \boldsymbol{\theta}_1), (\boldsymbol{x}_2, \boldsymbol{\theta}_2), \ldots, (\boldsymbol{x}_N, \boldsymbol{\theta}_N)$ from the joint posterior $\pi(\boldsymbol{x}, \boldsymbol{\theta}|\boldsymbol{y})$.

- Can be computed directly from MCMC output

- Asymptotically consistent

- Exhibit large or even infinite variance

---

[2]Newton M.A. and Raftery A.E. (1994) Approximate Bayesian inference with the weighted likelihood bootstrap *J. R. Stat. Soc. Ser. B. Stat. Methodol*, **56**, 3–48

# CHIB'S METHODS[3]

- Based on the observation that

$$\pi(\boldsymbol{y}) = \frac{\pi(\boldsymbol{y}|\boldsymbol{x}, \boldsymbol{\theta})\,\pi(\boldsymbol{x}, \boldsymbol{\theta})}{\pi(\boldsymbol{x}, \boldsymbol{\theta}|\boldsymbol{y})}$$

  for fixed $\boldsymbol{\theta}^*, \boldsymbol{x}^*$ (high–density posterior point) the log marginal likelihood can be estimated by

  $$\log \widehat{\pi}_{\mathrm{Chib}}(\boldsymbol{y}) = \log \pi(\boldsymbol{y}|\boldsymbol{x}^*, \boldsymbol{\theta}^*) + \log \pi(\boldsymbol{x}^*, \boldsymbol{\theta}^*) - \log \widehat{\pi}(\boldsymbol{x}^*, \boldsymbol{\theta}^*|\boldsymbol{y})$$

  $\Longrightarrow$ is estimated by breaking the parameter vector into appropriate blocks

- Required a separate MCMC run to calculate each block

---

[3]Chib S. (1995) Marginal likelihood from the Gibbs output *J. Amer. Statist. Assoc*, **90**, 1313–1321. Chib S. and Jeliazkov I. (2001) Marginal likelihood from the MH output *J. Amer. Statist. Assoc*, **96**, 270–281

MOTIVATION
ooo

METHODS
ooooooooo●oo

SIMULATION STUDIES
oooooo

REAL DATA ANALYSIS
ooooooooooo

CONCLUSIONS
oo

# POWER POSTERIORS[4]

- Power Posterior defined as

$$\pi(\boldsymbol{x}, \boldsymbol{\theta}|\boldsymbol{y}, t) \propto \pi(\boldsymbol{y}|\boldsymbol{x}, \boldsymbol{\theta})^t \, \pi(\boldsymbol{x}, \boldsymbol{\theta})$$

where $t \in [0, 1]$ is a temperature parameter

- The log of the marginal likelihood can be represented by

$$\log \pi(\boldsymbol{y}) = \int_0^1 \mathrm{E}_{\boldsymbol{x}, \boldsymbol{\theta}|\boldsymbol{y}, t} \Big\{ \log \pi(\boldsymbol{y}|\boldsymbol{x}, \boldsymbol{\theta}) \Big\} \, dt$$

$\implies$ is calculated numerically by discretising
$0 = t_0 < t_1 < \cdots < t_n = 1$, and then using trapezium rule.

---

[4]Friel N. and Pettitt A. N. (2008) Marginal likelihood estimation via power posteriors *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **70**, 589–607

# POWER POSTERIORS: EXAMPLE



- Obtain samples from the power posterior at each temperature $t_i$
- Variability depends
  - ➤ Number of $t_i$'s
  - ➤ Spacing of $t_i$'s
  - ➤ Number of MCMC samples
- Large number $\Longrightarrow$ more computational effort

MOTIVATION
○○○

METHODS
○○○○○○○○○○●

SIMULATION STUDIES
○○○○○○

REAL DATA ANALYSIS
○○○○○○○○○○

CONCLUSIONS
○○

# REVERSIBLE JUMP MCMC

# OUTLINE

1. **MOTIVATION**

2. **METHODS**

3. **SIMULATION STUDIES**

4. **REAL DATA ANALYSIS**

5. **CONCLUSIONS**

MOTIVATION
○○○

METHODS
○○○○○○○○○○○

SIMULATION STUDIES
●○○○○○

REAL DATA ANALYSIS
○○○○○○○○○○

CONCLUSIONS
○○

# SIMULATION STUDY: PNEMONOCOCCAL CARRIAGE[5]

- Household based longitudinal study on carriage of Streptococcus Pneumoniae

- Diagnostic tests obtained every 4 weeks
  - ➤ 10 months period
  - ➤ Classified as Negative / Positive

- The population is divided into two age groups:
  - ➤ Children 👧👦: under 5 years old
  - ➤ Adults 🧑🧑 : over 5 years old

---

[5]Touloupou et al. (2016) Model comparison with missing data using MCMC and importance sampling. arXiv 1512.04743

# MODEL DETAILS[6]

- Discrete time **S**usceptible-**I**nfected-**S**usceptible model

- The transition probabilities age group $i$ dependent:

$$P_i(S \longrightarrow I)_{\delta_t} = 1 - \exp\left\{ -\left( k_i + \frac{\beta_{Ci} I_C(t) + \beta_{Ai} I_A(t)}{(z-1)^w} \right) \cdot \delta t \right\}$$

$$P_i(I \longrightarrow S)_{\delta_t} = 1 - \exp\left( -\mu_i \cdot \delta t \right)$$

- 



[6]Melegaro et al. (2004) Estimating the transmission parameters of pneumococcal carriage in households. *Epidemiology and Infection*, **132**,

# RESULTS: MARGINAL LIKELIHOOD ESTIMATION



- $IS_{N_j} : N(\boldsymbol{\mu}, j\boldsymbol{\Sigma})$ • $IS_{t_d} : t_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ • $IS_{\mathrm{mix}} : 0.95 \times N(\boldsymbol{\theta}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) + 0.05\pi(\boldsymbol{\theta})$ • $\boldsymbol{\mu}, \boldsymbol{\Sigma}$: from MCMC
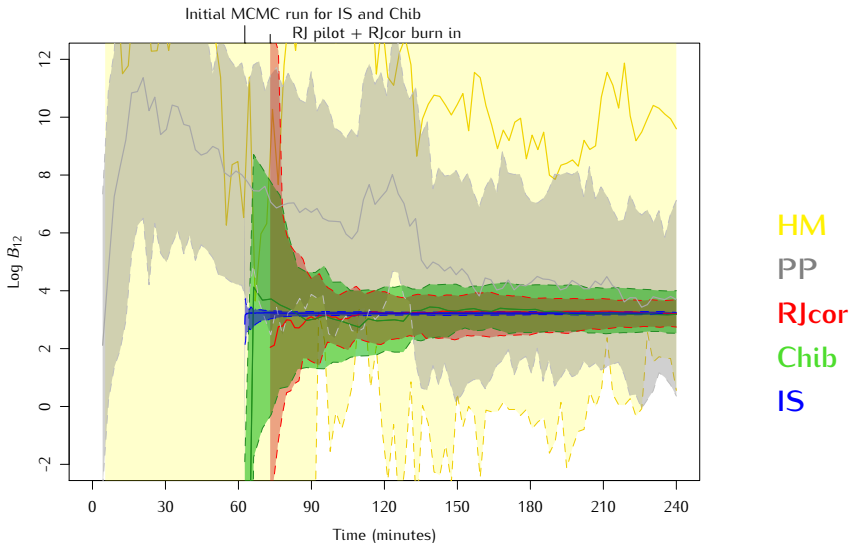
# HETEROGENEITY IN COMMUNITY ACQUISITION RATES

**Do adults and children acquire infection at the same rate?**

- We compare two models:
  - ➤ $\mathcal{M}_1 : k_A \neq k_C$
  - ➤ $\mathcal{M}_2 : k_A = k_C$

MOTIVATION
000

METHODS
00000000000

SIMULATION STUDIES
000000

REAL DATA ANALYSIS
0000000000

CONCLUSIONS
00

# RESULTS: BAYES FACTOR ESTIMATION



(a) Data simulated from model $\mathcal{M}_1$

(b) Data simulated from model $\mathcal{M}_2$

# RESULTS: EVOLUTION OF THE LOG BAYES FACTOR



HM
PP
**RJcor**
Chib
IS

## OUTLINE

1. **MOTIVATION**

2. **METHODS**

3. **SIMULATION STUDIES**

4. **REAL DATA ANALYSIS**

5. **CONCLUSIONS**

MOTIVATION
○○○

METHODS
○○○○○○○○○○○

SIMULATION STUDIES
○○○○○○

REAL DATA ANALYSIS
●○○○○○○○○○

CONCLUSIONS
○○

# STUDY DESIGNS

- Two longitudinal studies of *E. coli* O157:H7

|                   | Dataset 1     | Dataset 2  |
| ----------------- | ------------- | ---------- |
| Subjects          | 160 cattle    | 168 cattle |
| Study duration    | 14 weeks      | 22 weeks   |
| Sampling interval | 2 times/week  | 14 days    |

- Each sampling event included a
  - ➤ Faecal pat sample
  - ➤ Recto-anal mucosal swab (RAMS)

- Tests were assumed to have perfect specificity but imperfect sensitivity

# PATTERNS OF INFECTION



**Cattle in Pen 5**

# APPLICATION 1: E. COLI O157 IN FEEDLOT CATTLE

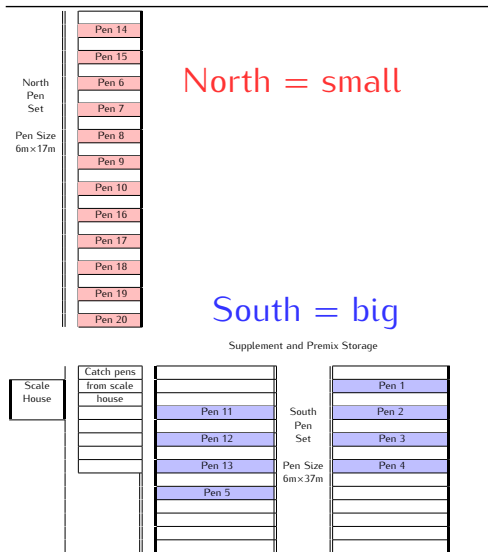**Do animals develop immunity over time?**

- We compare two models for infection period:
  - Geometric: lack of memory.
  - Negative Binomial: probability of recovery depends on duration of infection.

- The Negative Binomial is a generalisation of the Geometric:
  - Setting Negative Binomial dispersion parameter $\kappa = 1$ leads to Geometric.

MOTIVATION
000

METHODS
00000000000

SIMULATION STUDIES
000000

REAL DATA ANALYSIS
0000●000000

CONCLUSIONS
00

# APPLICATION 1: RESULTS



- **RJMCMC** and **IS** agree on the estimate of the Bayes factor

- **IS** estimator: faster convergence

- Bayes factor supports the Negative Binomial model

- The longer the colonization, the greater the probability of clearance – may indicate an immune response in the host
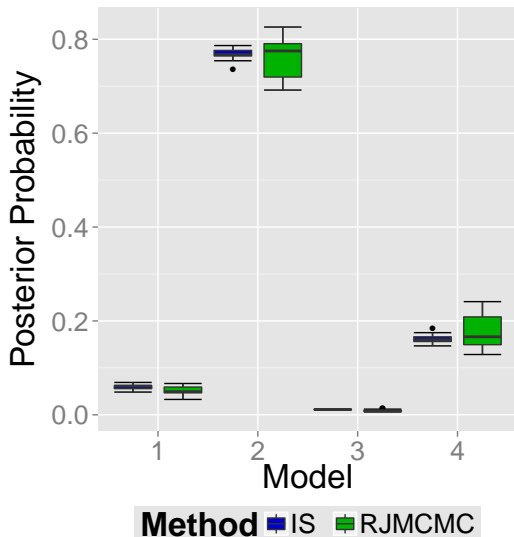
MOTIVATION
○○○

METHODS
○○○○○○○○○○○

SIMULATION STUDIES
○○○○○○

REAL DATA ANALYSIS
○○○○○●○○○○○

CONCLUSIONS
○○

# APPLICATION 2: ROLE OF PEN AREA/LOCATION

# APPLICATION 2: ROLE OF PEN AREA/LOCATION

> **Do north and south pens have different risk of infection?**

- Allow different external ($\alpha_s$, $\alpha_n$) and/or within-pen ($\beta_s$, $\beta_n$) transmission rates.

- Candidate models:

| Model | External North | External South | Within-pen North | Within-pen South |
|-------|--------|--------|--------|--------|
| 1 | $\alpha_n$ | $\alpha_s$ | $\beta_n$ | $\beta_s$ |
| 2 | $\alpha$ | $\alpha$ | $\beta_n$ | $\beta_s$ |
| 3 | $\alpha_n$ | $\alpha_s$ | $\beta$ | $\beta$ |
| 4 | $\alpha$ | $\alpha$ | $\beta$ | $\beta$ |

# APPLICATION 2: POSTERIOR PROBABILITIES



- **RJMCMC** and **IS** provide identical conclusions.

- Evidence to support different within-pen transmission rates.

- Animals in smaller pens more at risk of within-pen infection

# APPLICATION 3: INVESTIGATING TRANSMISSION BETWEEN PENS

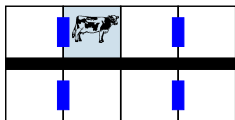Dataset 2: pens adjacent in a $12 \times 2$ rectangular grid.

- No direct contact across **feed buck**.
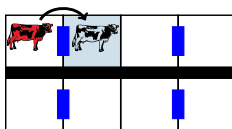- Shared **waterers** between pairs of adjacent pens.

MOTIVATION
OOO

METHODS
OOOOOOOOOOO

SIMULATION STUDIES
OOOOOO

REAL DATA ANALYSIS
OOOOOOOO●O

CONCLUSIONS
OO

# APPLICATION 3: INVESTIGATING TRANSMISSION BETWEEN PENS

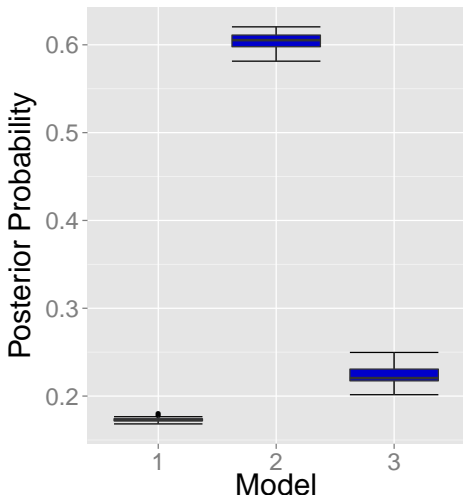## Do waterers spread infection?



(a) Model 1: No contacts between pens

(b) Model 2: Transmission via a waterer

(c) Model 3: Transmission via any boundary

# APPLICATION 3: POSTERIOR PROBABILITIES



- **RJMCMC**: hard to design efficient jump mechanism

- Using **IS** results still possible

- Evidence for transmission between pens sharing a waterer rather than another boundary

# OUTLINE

# CONCLUDING REMARKS

- Show how IS can be used to test epidemiological
  questions of interest

- In this study the importance sampling estimator
  outperformed existing tools
  ➣ Smallest Monte Carlo error

- Importance sampling approach very easy to implement and
  trivially parallelisable

- Bayes factors depend on choice of prior
  ➣ Simulations needed to avoid Lindley's paradox

# VARIATIONS/EXTENSIONS

- When the full conditional is not available we use a related full conditional
  - ➤ IS corrects for not using the true full conditional

- My collaborator Peter Neal used the particle filtering to estimate $\pi(x|\theta)$

- We recently applied Bridge Sampling for estimating the marginal likelihood
  - ➤ IS a special case
  - ➤ Slightly reduced variances
  - ➤ We use IS due to ease of implementation

MOTIVATION
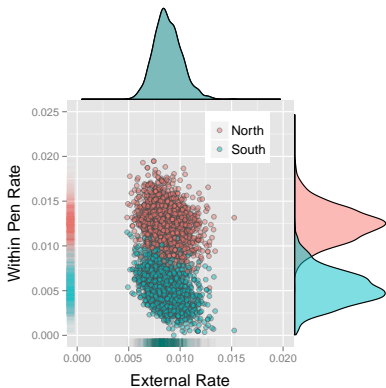○○○

METHODS
○○○○○○○○○○○

SIMULATION STUDIES
○○○○○○

REAL DATA ANALYSIS
○○○○○○○○○○

CONCLUSIONS
○○

# THANKS FOR LISTENING!

| Parameter | Symbol | Geometric | Negative Binomial |
|---|---|---|---|
| External transmission probability | $1 - e^{-\alpha}$ | 0.0090 [0.0064, 0.0117] | 0.0081 [0.0057, 0.0109] |
| Internal transmission probability | $1 - e^{-\beta}$ | 0.0107 [0.0077, 0.0141] | 0.0102 [0.0073, 0.0137] |
| Mean period of infection | $m$ | 8.9942 [7.7460, 10.4369] | 9.9740 [7.1977, 10.6487] |
| Shape parameter | $\kappa$ | — | 1.6245 [0.8361, 2.8972] |
| Initial probability of infection | $\mu$ | 0.1001 [0.0568, 0.1545] | 0.0997 [0.0557, 0.1546] |
| Sensitivity of RAJ test | $\theta_R$ | 0.7750 [0.7304, 0.8156] | 0.7771 [0.7311, 0.8203] |
| Sensitivity of faecal test | $\theta_F$ | 0.4639 [0.4206, 0.5073] | 0.4657 [0.4213, 0.5097] |

- Posterior mean of the parameters of each model along with the 95% credible interval in brackets.
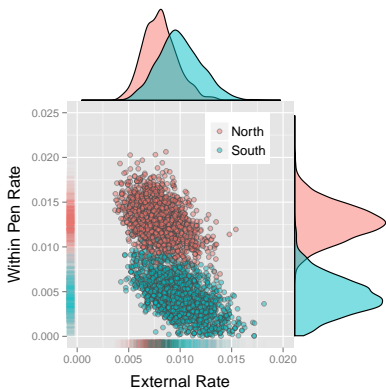
## PARAMETER ESTIMATION



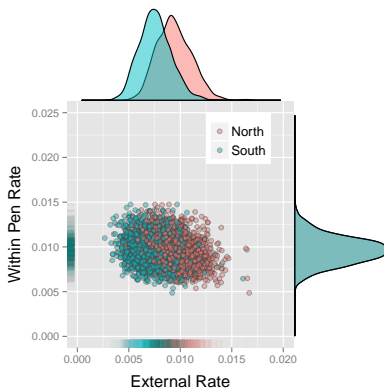(d) Model 2 – Posterior Prob 0.77          (e) Model 4 – Posterior Prob 0.16

# PARAMETER ESTIMATION



(f) Model 1 – Posterior Prob 0.06        (g) Model 3 – Posterior Prob 0.01

# THE CHOICE OF PRIOR MATTERS!
## SIMULATION STUDY: HETEROGENEITY IN TRANSMISSION RATES AMONG PENS