

Bayesian Emulation and Calibration of a Dynamic Epidemic Model for H1N1 Influenza

Marian Farah¹

Paul Birrell¹, Stefano Conti², Daniela De Angelis^{1,2}

¹MRC Biostatistics Unit, Cambridge, UK

²Health Protection Agency, London, UK

ICERM Bayesian Nonparametrics Workshop
September 19, 2012

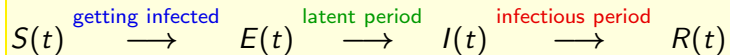
Motivation

- **Tracking** and **predicting** the behavior of an emerging epidemic is *essential* for a prompt public health response.
- **Inferential goals:**
 - **What is happening?** i.e., **real-time estimation** of the epidemic **parameters**.
 - **What is going to happen next?** i.e., **forecasting** the (short-term) **evolution** of the epidemic.
 - **What happened?** i.e., **“reconstructing”** the epidemic by **estimating** its **parameters** and **evolution dynamics**.
- **Noisy time-series** data coming from **different sources**.

- 1 Introduction: Epidemic modeling
- 2 Emulation and calibration of epidemic models
- 3 Preliminary results
- 4 Discussion

Introduction

- **Transmission model:**



- Transmission depends on the **virulence**, the **mixing patterns in the population**, and the **transition rates** among the S , E , I , and R states.
- Transmission **dynamics** are typically described by a system of **differential equations**.

Birrell et al. (2011) H1N1 model

$$S(t) \longrightarrow E(t) \longrightarrow I(t) \longrightarrow R(t)$$

Outline

Introduction

Methods

Results

Discussion

Birrell et al. (2011) H1N1 model

$$S(t) \longrightarrow E(t) \longrightarrow I(t) \longrightarrow R(t)$$

↓ incubation

Expected # of symptomatic individuals

↓ propensity to consult doctor

Expected # of doctor consultations

↓ delay in reporting

Expected # of reported cases, $\mu(\eta, t)$

Birrell et al. (2011) H1N1 model

$$S(t) \longrightarrow E(t) \longrightarrow I(t) \longrightarrow R(t)$$

↓ incubation

Expected # of symptomatic individuals

↓ propensity to consult doctor

Expected # of doctor consultations

↓ delay in reporting

Expected # of reported cases, $\mu(\eta, t)$

- $\eta = (\eta_1, \dots, \eta_6)$: underlying **parameters** of the epidemic.
- Proportion of symptomatic cases, propensity to consult, exponential growth rate, expected infectious period, a measure of the initial number of infected individuals, population interaction parameters.

Computational challenge

Marian Farah
Biostatistics
Cambridge

Outline

Introduction

Methods

Results

Discussion

- The **likelihood** of reported data, $z(t)$, $t = 1, \dots, T$, depends on μ .

- $$p(\eta \mid z_{\{1:T\}}, \mu) \propto \prod_{t=1}^T p(z(t); \mu(\eta, t)) \times p(\eta)$$
- $\mu(\eta, t)$ must be computed at **every MCMC** iteration.
- μ is **computationally expensive**.

Computational challenge

Marian Farah
Biostatistics
Cambridge

Outline

Introduction

Methods

Results

Discussion

- The **likelihood** of reported data, $z(t)$, $t = 1, \dots, T$, depends on μ .

- $$p(\eta \mid z_{\{1:T\}}, \mu) \propto \prod_{t=1}^T p(z(t); \mu(\eta, t)) \times p(\eta)$$
- $\mu(\eta, t)$ must be computed at **every MCMC** iteration.
- μ is **computationally expensive**.
- What about an efficient estimate?

Computer simulator

specify inputs $\eta = X$

$$X = \begin{pmatrix} x_{1,1} & \dots & x_{1,6} \\ x_{2,1} & \dots & x_{2,6} \\ \vdots & \vdots & \vdots \\ x_{n,1} & \dots & x_{n,6} \end{pmatrix}$$



run code

Birrell
et al.
(2011)



outputs

$$\begin{aligned} &\mu(\mathbf{x}_1, 1), \dots, \mu(\mathbf{x}_1, T) \\ &\mu(\mathbf{x}_2, 1), \dots, \mu(\mathbf{x}_2, T) \\ &\quad \vdots \\ &\mu(\mathbf{x}_n, 1), \dots, \mu(\mathbf{x}_n, T) \end{aligned}$$

Computer simulator

specify inputs $\eta = X$

$$X = \begin{pmatrix} x_{1,1} & \dots & x_{1,6} \\ x_{2,1} & \dots & x_{2,6} \\ \vdots & \vdots & \vdots \\ x_{n,1} & \dots & x_{n,6} \end{pmatrix}$$



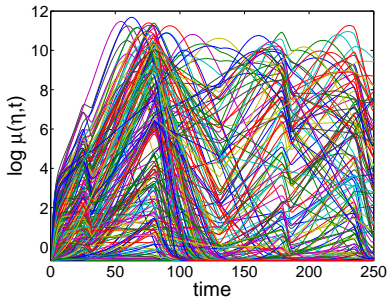
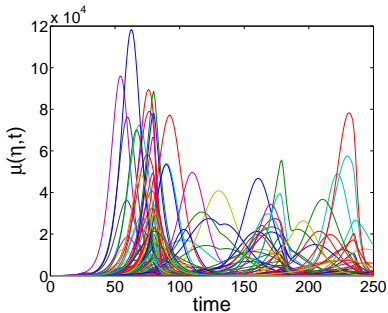
run code

Birrell
et al.
(2011)



outputs

$$\begin{aligned} &\mu(\mathbf{x}_1, 1), \dots, \mu(\mathbf{x}_1, T) \\ &\mu(\mathbf{x}_2, 1), \dots, \mu(\mathbf{x}_2, T) \\ &\vdots \\ &\mu(\mathbf{x}_n, 1), \dots, \mu(\mathbf{x}_n, T) \end{aligned}$$



Calibration and Emulation

Marian Farah
Biostatistics
Cambridge

Outline

Introduction

Methods

Results

Discussion

- **Calibration:** (e.g., Higdon et al., 2004)

Posterior inference for η through the simulator, μ , and “field” observed data $z(t)$,

$$\text{Observed} = \text{Reality} + \text{Error}$$

$$\text{Observed} = \text{Simulator} + \text{bias} + \text{Error}$$

z^\uparrow

μ^\uparrow

b^\uparrow

- $$p(\eta, b \mid z_{\{1:T\}}, \mu) \propto \prod_{t=1}^T p(z(t); \mu(\eta, t) + b) \times p(\eta)p(b)$$

- **Emulation:** (e.g., Kennedy and O’Hagan, 2000)

Estimating a slow computer *simulator* output, μ , using *fast* statistical model (an *emulator*), say $\hat{\mu}$.

Calibration and Emulation

Outline

Introduction

Methods

Results

Discussion

- **Idea:** (e.g., Bayarri et al., 2007a)
Replace the *slow* simulator output, μ , with the *fast* emulator estimation, $\hat{\mu}$, and obtain posterior inference for η through

- $$p(\eta, b \mid z_{\{1:T\}}, \hat{\mu}) \propto \prod_{t=1}^T p(z(t); \hat{\mu}(\eta, t) + b) \times p(\eta)p(b)$$

Emulation and calibration of dynamic models

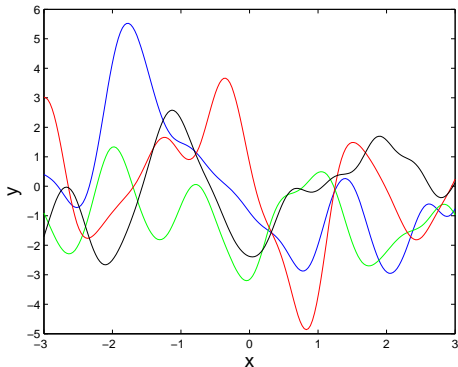
- A **deterministic** computer simulator is a **function** $f(\cdot)$ that maps **input** x to a unique **output** $y = f(x)$.
- The function $f(\cdot)$ is treated as **unknown** and given a **prior**.
- **Likelihood**: data are **runs of the simulator**, given a **design** over the **input space**, e.g., Latin Hypercube.
- **Emulator**: the **posterior** (predictive) distribution of $f(\cdot)$.

The Gaussian process

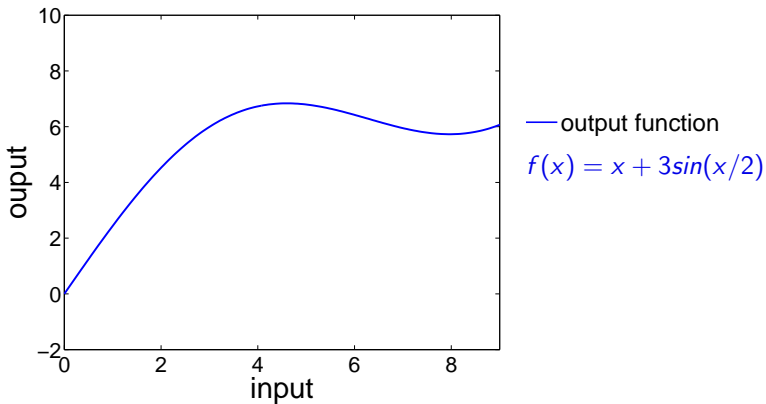
Marian Farah
Biostatistics
Cambridge

$$y(x) \sim GP(m(x), v c(x, x'))$$

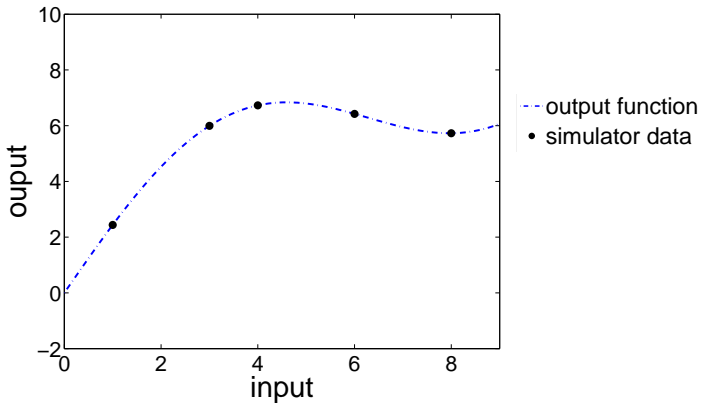
$m(\cdot)$, v , and $c(\cdot, \cdot)$ are the mean, variance, & correlation function (e.g., Neal 1998; Rasmussen & Williams 2006).



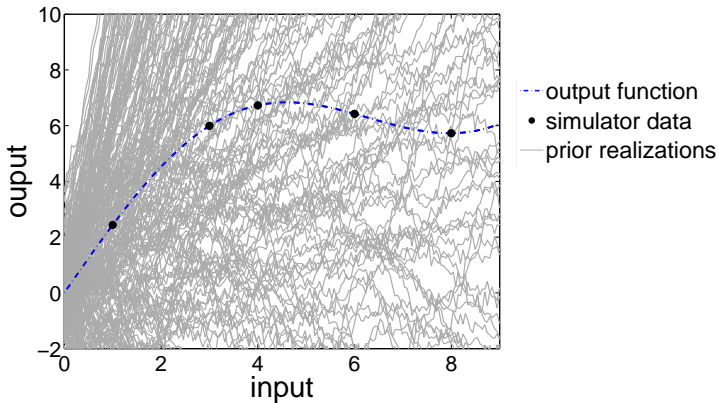
Toy example



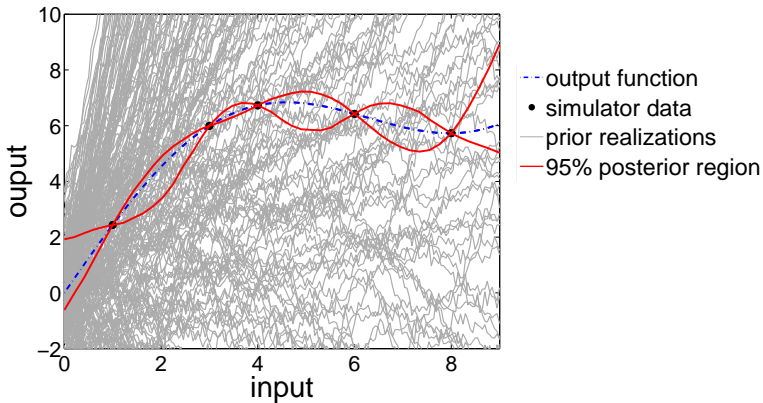
Toy example



Toy example

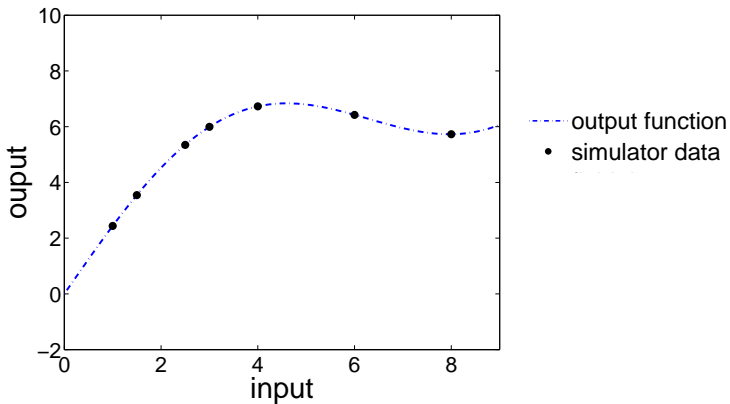


Toy example

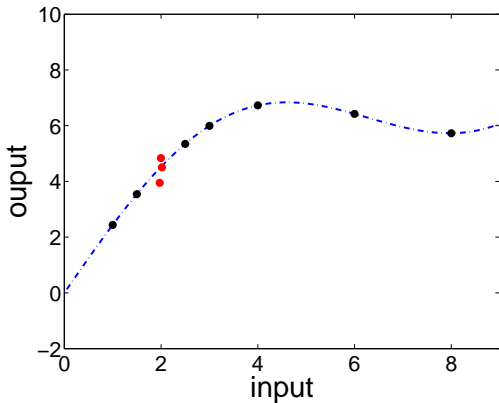


- Simulator: Specify $\mathbf{x} \rightarrow f(\mathbf{x})$.
- For $\mathbf{x} = \boldsymbol{\eta}$, $f(\boldsymbol{\eta})$ simulates a physical system.
- $\boldsymbol{\eta}$ is **uncertain**.
- **Calibration**: solving the **inverse-problem**, i.e., $\boldsymbol{\eta} \mid \mathbf{z}, f(\cdot)$.
- If $f(\cdot)$ is computationally expensive, it is emulated.
 - **Priors** for $\boldsymbol{\eta}$ and $f(\cdot)$.
 - **Likelihood**: data come from **field** observations and **simulator** runs.

Toy example



Toy example



--- output function

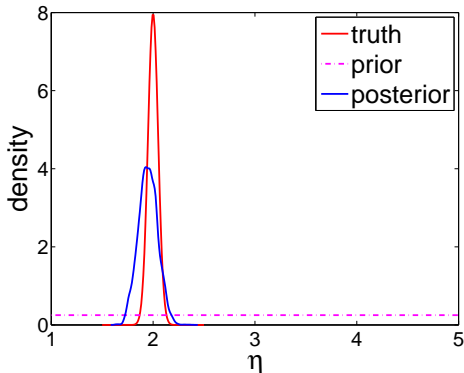
• simulator data

• field data

$$z \sim N(f(\eta), \sigma^2 = 0.3^2)$$

$$\eta \sim N(2, 0.05^2)$$

Toy example



- Assuming σ^2 is known.

- $y_t(x_i) = f(x_i, t)$ is the simulator output at input point x_i and time t .

$$\begin{array}{lcl} x_1 & \longrightarrow & y_1(x_1), y_2(x_1), \dots, y_T(x_1) \\ x_2 & \longrightarrow & y_1(x_2), y_2(x_2), \dots, y_T(x_2) \\ \vdots & & \vdots \\ x_n & \longrightarrow & y_1(x_n), y_2(x_n), \dots, y_T(x_n) \end{array}$$

- $y_t(x_i) = f(x_i, t)$ is the simulator output at input point x_i and time t .

$$\begin{array}{ccccccc} x_1 & \longrightarrow & y_1(x_1), & y_2(x_1), & \dots, & y_T(x_1) \\ x_2 & \longrightarrow & y_1(x_2), & y_2(x_2), & \dots, & y_T(x_2) \\ \vdots & & \vdots & \vdots & & \vdots \\ x_n & \longrightarrow & y_1(x_n), & y_2(x_n), & \dots, & y_T(x_n) \end{array}$$

- Need to model **three** types of **interdependencies**:
 - ① over the **input space**.
 - ② over time **within** each **time series**.
 - ③ **across series** of different input points.

Dynamic emulation

- **Modeling dependence over the input space alone**

Typically using a **Gaussian process** prior for outputs.

$$y(x) \sim GP(m(x), v c(x, x'))$$

- **Modeling dependence over the input space alone**

Typically using a **Gaussian process** prior for outputs.

$$y(x) \sim GP(m(x), v c(x, x'))$$

- **Modeling dependence for a single time series**

Typically, **TVAR(p)** model is used; e.g., $\mathbf{p} = 1$,

$$y_t(x) = \phi_t y_{t-1}(x) + \epsilon_t(x), \quad \epsilon_t(x) \sim N(0, v_t),$$

$$\phi_t = \phi_{t-1} + \omega_t, \quad \omega_t \sim N(0, w_t).$$

- **Linking across time series for different inputs** using a multivariate TVAR(p) model (Liu and West, 2009),

$$\begin{pmatrix} y_t(x_1) \\ \vdots \\ y_t(x_n) \end{pmatrix} = \begin{pmatrix} y_{t-1}(x_1) & \cdots & y_{t-p}(x_1) \\ \vdots & \ddots & \vdots \\ y_{t-1}(x_n) & \cdots & y_{t-p}(x_n) \end{pmatrix} \begin{pmatrix} \phi_{1,t} \\ \vdots \\ \phi_{p,t} \end{pmatrix} + \begin{pmatrix} \epsilon_t(x_1) \\ \vdots \\ \epsilon_t(x_n) \end{pmatrix}$$

- **Linking across time series for different inputs** using a multivariate TVAR(p) model (Liu and West, 2009),

$$\begin{pmatrix} y_t(x_1) \\ \vdots \\ y_t(x_n) \end{pmatrix} = \begin{pmatrix} y_{t-1}(x_1) & \cdots & y_{t-p}(x_1) \\ \vdots & \ddots & \vdots \\ y_{t-1}(x_n) & \cdots & y_{t-p}(x_n) \end{pmatrix} \begin{pmatrix} \phi_{1,t} \\ \vdots \\ \phi_{p,t} \end{pmatrix} + \begin{pmatrix} \epsilon_t(x_1) \\ \vdots \\ \epsilon_t(x_n) \end{pmatrix}$$

$$\text{Cov}(\epsilon_t(x_i), \epsilon_t(x_j)) = v_t c(x_i, x_j)$$

- $c(x_i, x_j)$ is the (i, j) element in the $n \times n$ **correlation** matrix induced by the **Gaussian process**.

- **Linking across time series for different inputs** using a multivariate TVAR(p) model (Liu and West, 2009),

$$\begin{pmatrix} y_t(x_1) \\ \vdots \\ y_t(x_n) \end{pmatrix} = \begin{pmatrix} y_{t-1}(x_1) & \cdots & y_{t-p}(x_1) \\ \vdots & \ddots & \vdots \\ y_{t-1}(x_n) & \cdots & y_{t-p}(x_n) \end{pmatrix} \begin{pmatrix} \phi_{1,t} \\ \vdots \\ \phi_{p,t} \end{pmatrix} + \begin{pmatrix} \epsilon_t(x_1) \\ \vdots \\ \epsilon_t(x_n) \end{pmatrix}$$

$$\text{Cov}(\epsilon_t(x_i), \epsilon_t(x_j)) = v_t c(x_i, x_j)$$

- $c(x_i, x_j)$ is the (i, j) element in the $n \times n$ **correlation matrix** induced by the **Gaussian process**.
- $\phi_t = \phi_{t-1} + \omega_t$, where $\phi_t = (\phi_{1t}, \dots, \phi_{pt})'$.

Birrell et al. (2011) simulator

Marian Farah
Biostatistics
Cambridge

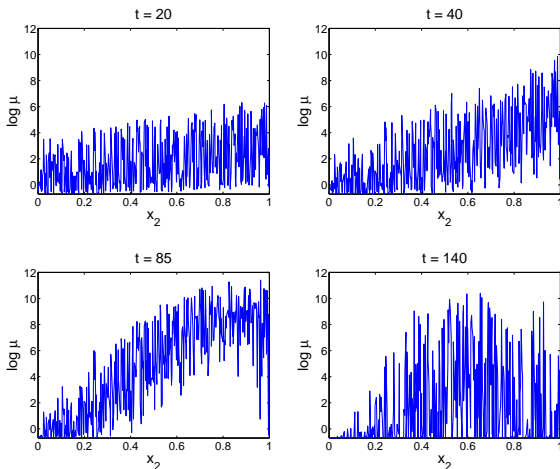
Outline

Introduction

Methods

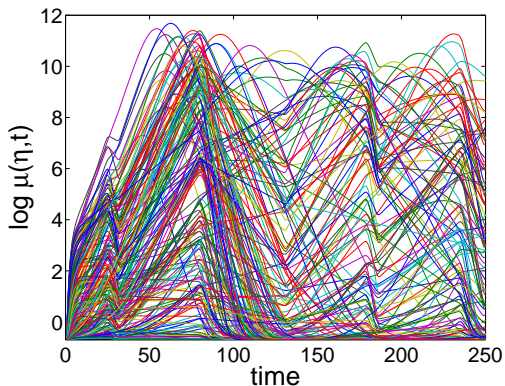
Results

Discussion



- x_2 : Exponential growth rate.

Birrell et al. (2011) simulator



- **Extending Liu and West (2009)**

- Modeling **input-dependent trends**:

$$y_t(x) = \phi_t y_{t-1}(x) + h(x)\beta_t + \epsilon_t$$

- Modeling **systematic temporal trend**:

$$y_t(x) = \theta_t + \phi_t y_{t-1}(x) + h(x)\beta_t + \epsilon_t$$

$$\begin{pmatrix} \theta_t \\ \phi_t \\ \beta_t \end{pmatrix} = \begin{pmatrix} \theta_{t-1} \\ \phi_{t-1} \\ \beta_{t-1} \end{pmatrix} + \begin{pmatrix} \omega_{1t} \\ \omega_{2t} \\ \omega_{3t} \end{pmatrix}$$

- Posterior inference through Forward-Filtering Backward-Sampling.

- **Two** sources of data:
 - **Simulator** data: $D^S = \{(y_t, \mathbf{x}); t = 1, \dots, T\}$. Model parameters are *specified* as inputs \mathbf{x} .
 - **“Field”** observed epidemic data $D^F = \{z_t; t = 1, \dots, T\}$. Model parameters, η , are unknown.
- **Two-stage** calibration (e.g., Bayarri et al., 2007b)
 - **Stage 1**: Estimate the emulator model parameters using only D^S .
 - **Stage 2**: Model z_t using a parametric distribution centered on the emulator model. Then, conditional on stage 1, estimate $p(\eta | D^F, D^S)$.

Results

Validating the emulator

Marian Farah
Biostatistics
Cambridge

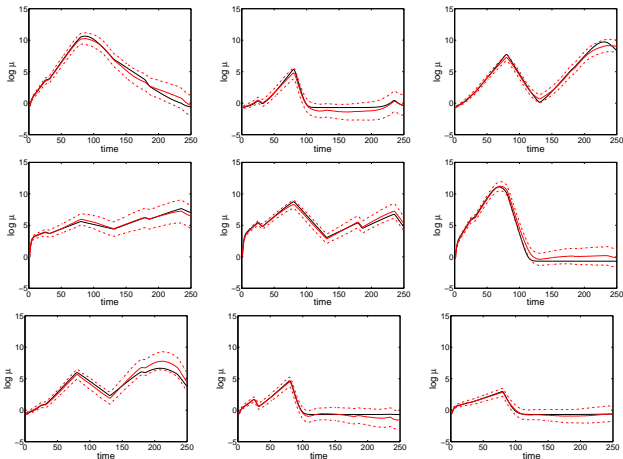
Outline

Introduction

Methods

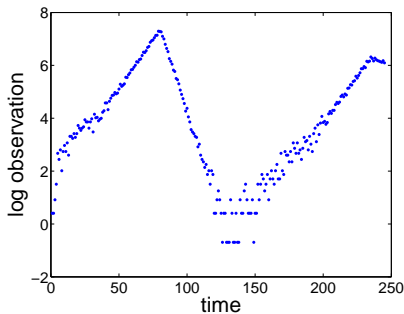
Results

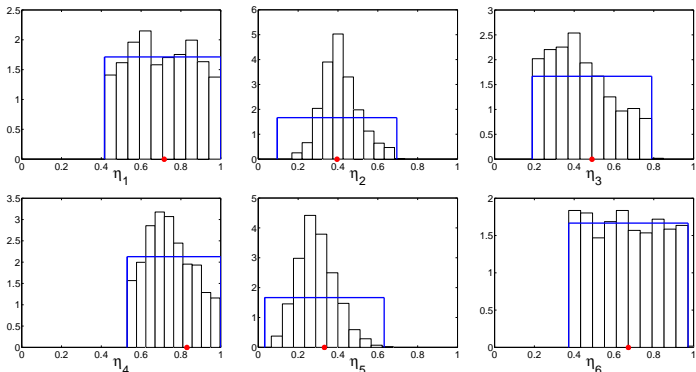
Discussion



- Simulation runs (black), emulator's median & 95% region (red).
- Plots based on a MVTVAR(1) and Gaussian correlation function.

- Generated synthetic epidemic data.
- Set $\eta = \eta_0$. Then, $z \sim \text{Poisson}(\mu(\eta_0, t))$





• Truth — Prior

- η_2 is exponential growth rate, and η_5 is effect of summer holiday on population interaction.

Discussion

- What we have done:
 - **Estimation** of epidemic dynamics by combining a **statistical emulator** with **reported epidemic data**.
 - **Dynamic emulation** through modeling dependencies across time and epidemic parameter space.
- Still to do:
 - Consider different **age** groups in the population.
 - Incorporate **additional sources** of information.
 - Real-time calibration and forecasting using epidemic data.
 - ...

- What we have done:
 - **Estimation** of epidemic dynamics by combining a **statistical emulator** with **reported epidemic data**.
 - **Dynamic emulation** through modeling dependencies across time and epidemic parameter space.
- Still to do:
 - Consider different **age** groups in the population.
 - Incorporate **additional sources** of information.
 - Real-time calibration and forecasting using epidemic data.
 - ...

Thank you!

- Bayarri, M., Berger, J., Paulo, R., Sacks, J., Cafeo, J., Cavendish, J., Lin, C., and Tu, J. (2007a), "A framework for validation of computer models," *Technometrics*, 49, 138–154.
- Bayarri, M. J., Berger, J. O., Cafeo, J., Garcia-Donato, G., Liu, F., Palomo, J., Parthasarathy, R., Paulo, R., Sacks, J., and Walsh, D. (2007b), "Computer model validation with functional output," *Annals of Statistics*, 35, 1874–1906.
- Birrell, P. J., Ketsetzis, G., Gay, N. J., Cooper, B. S., Presanis, A. M., Harris, R. J., Charlett, A., Zhang, X.-S., White, P. J., Pebody, R. G., and De Angelis, D. (2011), "Bayesian modeling to unmask and predict influenza A/H1N1pdm dynamics in London," *Proceedings of the National Academy of Sciences*.
- Kennedy, M. C. and O'Hagan, A. (2000), "Predicting the output from a complex computer code when fast approximations are available," *Biometrika*, 87, 1–13.
- Liu, F. and West, M. (2009), "A Dynamic Modelling Strategy for Bayesian Computer Model Emulation," *Bayesian Analysis*, 4, 393–412.