

# Statistical Inference for epidemic models approximated by diffusion processes

Romain GUY<sup>1,2</sup>

Joint work with C. Larédo<sup>1,2</sup> and E. Vergu<sup>1</sup>

<sup>1</sup> UR 341, MIA, INRA, Jouy-en-Josas

<sup>2</sup> UMR 7599, LPMA, Université Paris Diderot

ANR MANEGE, Université Paris 13

30 janvier 2013

## Outline

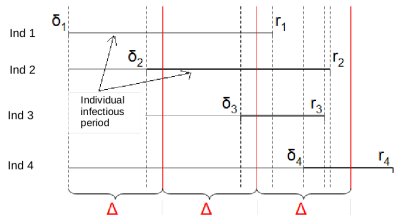
Provide a framework for estimating key parameters of epidemics

- 1 Characteristics of the epidemic process
  - Constraints imposed by the observation of the epidemic process
  - Simple mechanistic models
- 2 Various mathematical approaches for epidemic spread
  - Natural approach: Markov jump process
  - First approximation by ODEs
  - Gaussian approximation of the Markov jump process
  - Diffusion approximation of the Markov jump process
- 3 Inference for discrete observations of diffusion or Gaussian processes with small diffusion coefficient
  - Contrast processes for fixed or large number of observations
  - Correction of a non asymptotic bias
  - Comparison of estimators on simulated epidemics
- 4 Epidemics incompletely observed: partially and integrated diffusion processes (Work in progress)
  - Back to epidemic data
  - Inference approach: Work in progress

## Outline

- 1 Characteristics of the epidemic process
  - Constraints imposed by the observation of the epidemic process
  - Simple mechanistic models
- 2 Various mathematical approaches for epidemic spread
  - Natural approach: Markov jump process
  - First approximation by ODEs
  - Gaussian approximation of the Markov jump process
  - Diffusion approximation of the Markov jump process
- 3 Inference for discrete observations of diffusion or Gaussian processes with small diffusion coefficient
  - Contrast processes for fixed or large number of observations
  - Correction of a non asymptotic bias
  - Comparison of estimators on simulated epidemics
- 4 Epidemics incompletely observed: partially and integrated diffusion processes (Work in progress)
  - Back to epidemic data
  - Inference approach: Work in progress

## Incomplete Data

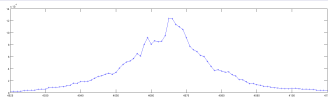


Individual Tracking	{ Ind 1, Ind 2 }	{ Ind 1, Ind 2, Ind 3 }	{ Ind 4 }
Number of Infected	2	3	1
Incidence (new infected)	2	1	1

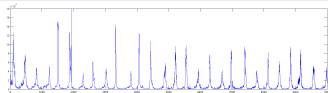
## Different dynamics framework

Ex.: Influenza like illness cases (Sentinelles surveillance network)

One outbreak study



Recurrent outbreaks study



## Imperfect data

- Incomplete observations
- Temporally aggregated
- Sampling & reporting error
- Unobserved cases

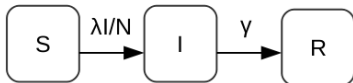
Main goal: key parameter estimation

- Basic reproduction number,  $R_0$  (nb. of secondary cases generated by one primary case in an entirely susceptible population)
- Average infectious time period ( $d$ )

## Compartmental representation of the population dynamics

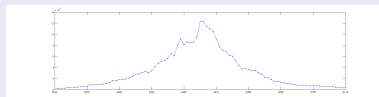
Define: Nb. of health states, possible transitions and associated rates Notations  $N$  : population size,  $\lambda$  : transmission rate,  $\gamma$  : recovery rate  
 $S, I, R$  : numbers of susceptible, infected, removed individuals

## One of the simplest model: SIR



Closed population  $\Rightarrow N = S + I + R$   
 Well-mixing population  
 $\Rightarrow (S, I) \xrightarrow{\lambda SI/N} (S-1, I+1)$

Convenient to study one epidemics



Key parameters:  $R_0 = \frac{\lambda}{\gamma}$ ,  $d = \frac{1}{\gamma}$

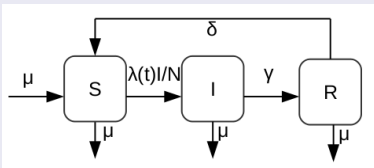
Summary: coefficients  $\alpha_L$ 

$(S, I) \rightarrow (S-1, I+1) = (S, I) + (-1, 1)$  at rate  $\alpha_{(-1,1)}(S, I) = \lambda S \frac{I}{N}$  and  
 $(S, I) \rightarrow (S, I-1) = (S, I) + (0, -1)$  at rate  $\alpha_{(0,-1)}(S, I) = \gamma I$

Natural extensions of the SIR model

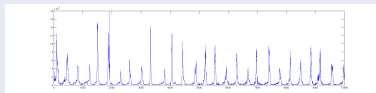
- Increase the number of health states : e.g. Exposed Class  $\Rightarrow$  SEIR model
- Additional transitions: e.g.  $(S, I) \rightarrow (S - 1, I)$  (vaccination)

Temporal dependence: SIRS with seasonality in transmission and demography



$\delta$ : waning immunity rate (years)  
 $\mu$ : demographic renewal rate (decades)  
 $\lambda(t) = \lambda_0(1 + \lambda_1 \sin(2\pi \frac{t}{T_{per}}))$   
 $\lambda_1 = 0 \Rightarrow$  oscillations vanishes

Suited to study recurrent epidemics



Key parameters:  $R_0^{Moy} = \frac{\lambda_0}{\gamma + \mu}$ ,  $d = \frac{1}{\gamma}$

Summary:

$$\alpha_{(-1,1)}(t, S, I) = \lambda(t)S \frac{I}{N}$$

$$\alpha_{(1,0)}(S, I) = N\mu + \delta(N - S - I),$$

$$\alpha_{(-1,0)}(S, I) = \mu S$$

$$\alpha_{(0,-1)}(S, I) = (\mu + \gamma)I$$

## Outline

- 1 Characteristics of the epidemic process
  - Constraints imposed by the observation of the epidemic process
  - Simple mechanistic models
- 2 Various mathematical approaches for epidemic spread
  - Natural approach: Markov jump process
  - First approximation by ODEs
  - Gaussian approximation of the Markov jump process
  - Diffusion approximation of the Markov jump process
- 3 Inference for discrete observations of diffusion or Gaussian processes with small diffusion coefficient
  - Contrast processes for fixed or large number of observations
  - Correction of a non asymptotic bias
  - Comparison of estimators on simulated epidemics
- 4 Epidemics incompletely observed: partially and integrated diffusion processes (Work in progress)
  - Back to epidemic data
  - Inference approach: Work in progress

## Markov jump process



$$\alpha_{(-1,1)}(S, I) = \lambda S \frac{I}{N}, \quad \alpha_{(0,-1)}(S, I) = \gamma I$$

Notations:

$$E = \{0, \dots, N\}^d$$

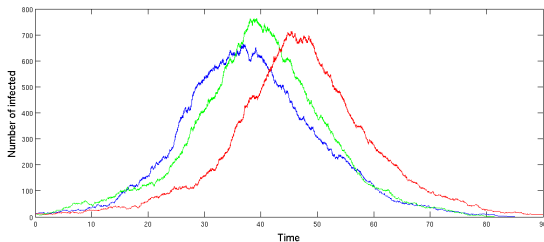
$\forall L \in E^- = \{-N, \dots, N\}^d$ , we define  $\alpha_L(\cdot) : E \rightarrow [0, +\infty[$

We define  $(Z_t)$  the Markov jump process on  $E$  with  $Q$ -matrix:  $q_{X,Y} = \alpha_{Y-X}(X)$

Assume  $\alpha(X) = \sum_{L \in E^-} \alpha_L(X) < +\infty \Rightarrow$  Sojourn time  $\text{Exp}(\alpha(X))$

## Easily simulated (using Gillespie algorithm)

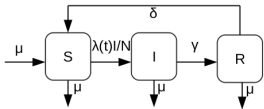
3 realizations for  $N = 10000$ ,  $\lambda = 0.5$ ,  $\gamma = 1/3$ ,  $(S_0, I_0) = (9990, 10)$





## Interest of the deterministic approach

$$\lambda(t) = \lambda_0(1 + \lambda_1 \sin(2\pi t / T_{per}))$$



## SIRS ODE solution

$$\frac{ds}{dt} = \mu(1 - s) + \delta(1 - s - i) - \lambda(t)si$$

$$\frac{di}{dt} = \lambda(t)si - (\mu + \gamma)i$$

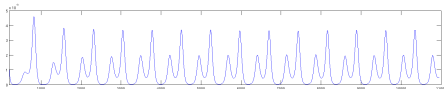
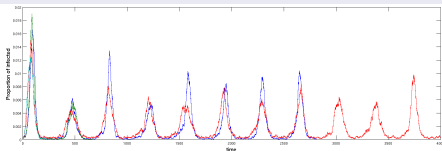
$$\lambda(t) = \lambda_0(1 + \lambda_1 \sin(2\pi t / T))$$

$$(s(0), i(0)) = \frac{Z_0}{N}$$

## ODE trajectory

Some trajectories of the SIRS Markov proc. :

$$N = 10^5, R_0 = 1.5, d = 3, \frac{1}{\delta T_{per}} = 2, \frac{1}{\mu T_{per}} = 50$$



## Drawbacks of the Markov jump approach

- $N = 10^7$  : more than  $10^5$  events in one week (MLE: observation of all the jumps required)
- Extinction probability non negligible

## Link between the two approaches

As  $N \rightarrow +\infty$  we have  $\frac{Z_t}{N} \xrightarrow[N \rightarrow \infty]{} x(t)$ , where

$x(t)$  is the deterministic solution of the ODE:

$$\frac{dx(t)}{dt} = b(x(t))$$

Function  $b$  is explicit

## Beyond deterministic limit : Gaussian process

Additional assumption: smooth version of  $\alpha_L, \beta_L$

We have  $\alpha_L : E \rightarrow (0, +\infty)$  transition rate :  $X \xrightarrow{\alpha_L(x)} X + l$

Assume  $\beta_L : [0, 1]^d \rightarrow [0, +\infty[$  well define and regular :

$$\forall x \in [0, 1]^d, \frac{1}{N} \alpha_L(\lfloor Nx \rfloor) \xrightarrow{N \rightarrow \infty} \beta_L(x)$$

$$\text{SIR: } \alpha_{(-1,1)}(S, I) = \lambda S \frac{I}{N} \Rightarrow \beta_{(-1,1)}(x) = \lambda x_1 x_2, \alpha_{(0,1)}(S, I) = \gamma I \Rightarrow \beta_{(0,-1)}(x) = \gamma x_2$$

Definition of function  $b \left( \frac{dx(t)}{dt} = b(x(t)) \right)$

$$b(x) = \sum_{L \in E^-} L \beta_L(x), \text{ SIR: } b((\lambda, \gamma), (s, i)) = \begin{pmatrix} -1 \\ 1 \end{pmatrix} \lambda si + \begin{pmatrix} 0 \\ 1 \end{pmatrix} \gamma i = \begin{pmatrix} -\lambda si \\ \lambda si + \gamma i \end{pmatrix}$$

ODE approximation : no longer dependance w.r.t.  $N$

Asymptotic expansion w.r.t.  $N$  : Gaussian process

$\sqrt{N} \left( \frac{Z_t}{N} - x(t) \right) \xrightarrow{N \rightarrow \infty} g(t)$  where  $g(t)$  centered Gaussian process:

$$dg(t) = \frac{\partial b}{\partial x}(x(t))g(t)dt + \sigma(x(t))dB_t, \text{ where } \sigma^t \sigma(x) = \Sigma(x) = \sum_{L \in E^-} L^t L \beta_L(x)$$

$$\text{SIR: } \Sigma((\lambda, \gamma), (s, i)) = \lambda si \begin{pmatrix} -1 \\ 1 \end{pmatrix} \begin{pmatrix} -1 & 1 \end{pmatrix} + \gamma i \begin{pmatrix} 0 \\ 1 \end{pmatrix} \begin{pmatrix} 0 & 1 \end{pmatrix} = \begin{pmatrix} \lambda si & -\lambda si \\ -\lambda si & \lambda si + \gamma i \end{pmatrix}$$

$$\text{Cholesky algorithm } \Rightarrow \sigma(x) = \begin{pmatrix} \sqrt{\lambda si} & 0 \\ -\sqrt{\lambda si} & \sqrt{\gamma i} \end{pmatrix}$$

## Infinitesimal generator approach: diffusion approximation

Renormalized version of the Markov jump process ( $\tilde{Z}_t$ ): $(\tilde{Z}_t)$  Markov jump process on  $E$  with transition rates  $q_{X,Y} = N\beta_{Y-X}(X)$  $\tilde{Z}_t = \frac{Z_t}{N}$  normalized processAsymptotic development of the infinitesimal generator of  $Z_t$  (Ethier & Kurtz (86))

Generator of  $\tilde{Z}_t$ :  $\mathcal{A}f(x) = \sum_{L \in E^-} NL\beta_L(x) (f(x+L) - f(x))$

$$\Rightarrow \text{Generator of } \tilde{Z}_t: \bar{\mathcal{A}}f(x) = b(x) \cdot \nabla f(x) + \frac{1}{N} \sum_{i,j=1}^d \Sigma_{i,j}(x) \frac{\partial^2 f}{\partial x_i \partial x_j}(x) + O\left(\frac{1}{N^2}\right)$$

Dropping negligible terms leads to the generator of a diffusion process:

$$dX_t = b(X_t)dt + \frac{1}{\sqrt{N}}\sigma(X_t)dB_t, \text{ where } b(x) = \sum_{L \in E^-} L\beta_L(x), \Sigma(x) = \sum_{L \in E^-} L^t L\beta_L(x)$$

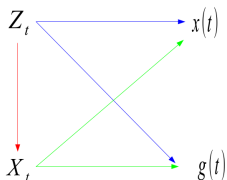
Temporal dependence  $\beta_L(t, x)$ : generator approach no longer availableDecomposition of the diffusion process using Gaussian process ( $\epsilon = \frac{1}{\sqrt{N}}$ )

Taylor's stochastic formula (Wentzell-Freidlin(79), Azencott (82))

Let  $dX_t = b(X_t)dt + \epsilon\sigma(X_t)dB_t, X_0 = x_0$ Then, under regularity assumptions,  $X_t = x(t) + \epsilon g(t) + O_{\mathbb{P}}(\epsilon^2)$

## Links between approximations

- $Z_t$  Markov jump process on  $E$  with transition  $q_{X,Y} = \alpha_{Y-X}(X)$
- $\bar{Z}_t$  Markov jump process on  $E/N$  with transition  $q_{x,y} = N\beta_{\lfloor Nx \rfloor - \lfloor Ny \rfloor}(x)$
- $\bar{Z}_t \sim x(t) + \frac{1}{\sqrt{N}}g(t)$  with  $x(\cdot)$  the ODE solution and  $g$  a Gaussian process
- $X_t$ :  $dX_t = b(x(t))dt + \frac{1}{\sqrt{N}}\sigma(x(t))dB_t$  diffusion with small diffusion coefficient  
and  $\mathbb{P}\left\{ \sup_{0 \leq t \leq T} \|\bar{Z}_t - X_t\| > C_T \frac{\log(N)}{N} \right\} \xrightarrow{N \rightarrow \infty} 0$
- $X_t = x(t) + \frac{1}{\sqrt{N}}g(t) + \mathcal{O}_{\mathbb{P}}\left(\frac{1}{N}\right)$



- Diffusion approximation
- Expansion in  $N$  of the process
- Taylor's stochastic expansion

Important : All mathematical representations completely defined by  $(\alpha_L)$

## Outline

- 1 Characteristics of the epidemic process
  - Constraints imposed by the observation of the epidemic process
  - Simple mechanistic models
- 2 Various mathematical approaches for epidemic spread
  - Natural approach: Markov jump process
  - First approximation by ODEs
  - Gaussian approximation of the Markov jump process
  - Diffusion approximation of the Markov jump process
- 3 Inference for discrete observations of diffusion or Gaussian processes with small diffusion coefficient
  - Contrast processes for fixed or large number of observations
  - Correction of a non asymptotic bias
  - Comparison of estimators on simulated epidemics
- 4 Epidemics incompletely observed: partially and integrated diffusion processes (Work in progress)
  - Back to epidemic data
  - Inference approach: Work in progress

## Classical Estimators

## Maximum likelihood estimation for the SIR Markov Jump process (Andersson &amp; Britton (00))

- Observation of all jumps
- Analytic expression of the estimators for SIR model:

$$\hat{\lambda} = N \frac{s(0) - s(T)}{\int_0^T S(t)I(t)dt}, \quad \hat{\gamma} = \frac{s(0) + i(0) - s(T) - i(T)}{\int_0^T I(t)dt}$$

- $\sqrt{N}(\hat{\theta}_{MLE} - \theta_0) \xrightarrow{N \rightarrow \infty} \mathcal{N}(0, I_b^{-1}(\theta_0))$ , where

$$I_b^{-1}((\lambda_0, \gamma_0)) = \begin{pmatrix} \frac{\lambda_0^2}{s(0) - s(T)} & 0 \\ 0 & \frac{\gamma_0^2}{s(0) + i(0) - s(T) - i(T)} \end{pmatrix}$$

## Maximum likelihood estimation for homoscedastic observations of the ODE

- $n$  observations at  $n$  discrete times  $t_k$  of  $x_\theta(t_k) + \xi_k$ , with  $\xi_k \sim \mathcal{N}(0, C_N(\theta_0)I_d)$
- MLE=LSE
- $\sqrt{n}(\hat{\theta}_{LSE} - \theta_0) \xrightarrow{n \rightarrow \infty} \mathcal{N}(0, I^N(\theta_0))$

## Specificities of the statistical framework

Model:

We define  $X_t$ :  $dX_t = b(\theta_1, X_t)dt + \epsilon\sigma(\theta_2, X_t)dB_t$ ,  $X_0 = x_0 \in \mathbb{R}^d$

⇒ separation of the parameters  $\theta_1, \theta_2$  required (not estimated at the same rate)

Continuous observation of the diffusion on  $[0, T]$  (Kutoyants (80))

$$\text{MLE} : \epsilon^{-1} \left( \theta_1^{MLE} - \theta_1^0 \right) \rightarrow \mathcal{N} \left( 0, I_b(\theta_1^0, \theta_2^0)^{-1} \right)$$

Existing discrete observation results : estimation of  $\theta_2$  at rate  $\sqrt{n}$

Observations:

We observe  $X_{t_k}$  for  $t_k = k\Delta$ ,  $k \in \{0, \dots, n\}$ ,  $t_k \in [0, T]$  ( $n\Delta = T$ ),  $T$  is fixed

Two different asymptotics:  $\epsilon \rightarrow 0$  &  $n$  ( $\Delta$ ) is fixed //  $\epsilon \rightarrow 0$  &  $n \rightarrow \infty$  ( $\Delta \rightarrow 0$ )

Notation:  $\eta = (\theta_1, \theta_2)$

SIR:  $\epsilon = \frac{1}{\sqrt{N}}$ ,  $\theta_1 = \theta_2 = \theta = (\lambda, \gamma)$ , and  $I_b(\theta_1^0, \theta_2^0)$  equals the Markov jump process Fisher Information matrix

Main idea: study of  $g_\eta(t)$  (Multidimensionnal generalization of Genon-Catalot(90))

Gaussian process:  $Y_t = x_{\theta_1}(t) + \epsilon g_\eta(t)$ ,  $n$  obs. at regular time intervals  $t_k = k\Delta$ , for  $k = 1, \dots, n$ .

Definition: Resolvent matrix of the linearized ODE system  $\Phi_{\theta_1}$

Let  $\Phi_{\theta_1}$  be the invertible matrix solution of

$$\frac{d\Phi_{\theta_1}}{dt}(t, t_0) = \frac{\partial b}{\partial x}(x_{\theta_1}(t))\Phi_{\theta_1}(t, t_0), \text{ with } \Phi_{\theta_1}(t_0, t_0) = I_d.$$

Important property of  $g_\eta$

$$g_\eta(t_k) = \Phi_{\theta_1}(t_k, t_{k-1})g_\eta(t_{k-1}) + \sqrt{\Delta}Z_k^\eta$$

$(Z_k^\eta)_{k \in \{1, \dots, n\}}$  independent Gaussian vectors with covariance matrix  $S_k^\eta$

$$S_k^\eta = \frac{1}{\Delta} \int_{t_{k-1}}^{t_k} \Phi_{\theta_1}(t_k, s) \Sigma(\theta_2, x_{\theta_1}(s)) {}^t \Phi_{\theta_1}(t_k, s) ds$$

Function of the observations

Let  $y \in \mathcal{C}([0, T], \mathbb{R}^d)$

$$N_k(\theta_1, y) = y(t_k) - x_{\theta_1}(t_k) - \Phi_{\theta_1}(t_k, t_{k-1}) [y(t_{k-1}) - x_{\theta_1}(t_{k-1})] (= \epsilon \sqrt{\Delta} Z_k^\eta)$$



Back to the diffusion process:  $dX_t = b(\theta_1, X_t)dt + \epsilon\sigma(\theta_2, X_t)dB_t$

$(Z_k^\eta)$  Gaussian family  $\Rightarrow$  Likelihood tractable:

$$-L_{\Delta, \epsilon}(\eta) = \epsilon^2 \sum_{k=1}^n \log[\det(S_k^\eta)] + \frac{1}{\Delta} \sum_{k=1}^n {}^t N_k(\theta_1, Y)(S_k^\eta)^{-1} N_k(\theta_1, Y)$$

$\hat{\theta}_2$  has good properties as  $\epsilon \rightarrow 0$ , only if  $\Delta \rightarrow 0$  ( $n \rightarrow +\infty$ )

1.  $n$  fixed,  $\epsilon \rightarrow 0$ : General case (low frequency contrast with  $\theta_2$  unknown)

$$\bar{U}_\epsilon(\theta_1) = \frac{1}{\Delta} \sum_{k=1}^n {}^t N_k(\theta_1, X) N_k(\theta_1, X) \Rightarrow \text{Associated MCE } \bar{\theta}_{1, \epsilon} = \underset{\theta_1 \in \Theta}{\operatorname{argmin}} \bar{U}_\epsilon(\theta_1)$$

2.  $n$  fixed,  $\epsilon \rightarrow 0$ : Case  $\theta_2 = f(\theta_1)$  (low frequency contrast with information on  $\theta_2$ )

$$\tilde{U}_\epsilon(\theta_1) = \frac{1}{\Delta} \sum_{k=1}^n {}^t N_k(\theta_1, X) (\tilde{S}_k^{\theta_1, f(\theta_1)})^{-1} N_k(\theta_1, X) \Rightarrow \tilde{\theta}_{1, \epsilon} = \underset{\theta_1 \in \Theta}{\operatorname{argmin}} \tilde{U}_\epsilon(\theta_1)$$

3.  $n \rightarrow \infty$ ,  $\epsilon \rightarrow 0$  (high frequency contrast)

$$\begin{aligned} \check{U}_{\Delta, \epsilon}(\theta_1, \theta_2) &= \epsilon^2 \sum_{k=1}^n \log[\det(\Sigma(\theta_2, X_{t_{k-1}}))] + \frac{1}{\Delta} \sum_{k=1}^n {}^t N_k(\theta_1, X) \Sigma^{-1}(\theta_2, X_{t_{k-1}}) N_k(\theta_1, X) \\ \Rightarrow \check{\theta}_{1, \epsilon, \Delta}, \check{\theta}_{2, \epsilon, \Delta} &= \underset{\eta \in \Theta}{\operatorname{argmin}} \check{U}_{\epsilon, \Delta}(\eta) \end{aligned}$$

What kind of distance is minimized: comparison with Least squares

$$N_k(\theta_1, y) = y(t_k) - x_{\theta_1}(t_k) - \Phi_{\theta_1}(t_k, t_{k-1}) [y(t_{k-1}) - x_{\theta_1}(t_{k-1})]$$

$N = 1000$ ,  $R_0 = 1.5$ ,  $d = 3$  days, 1 obs/day,  $T = 50$  days

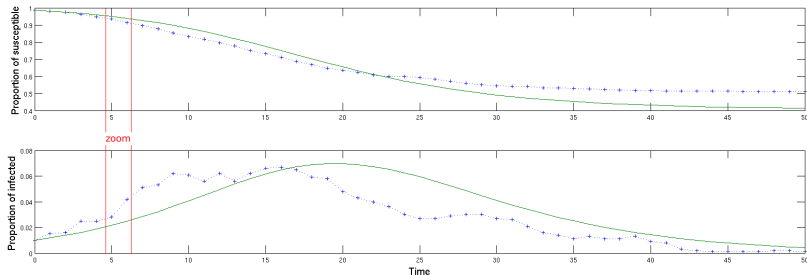
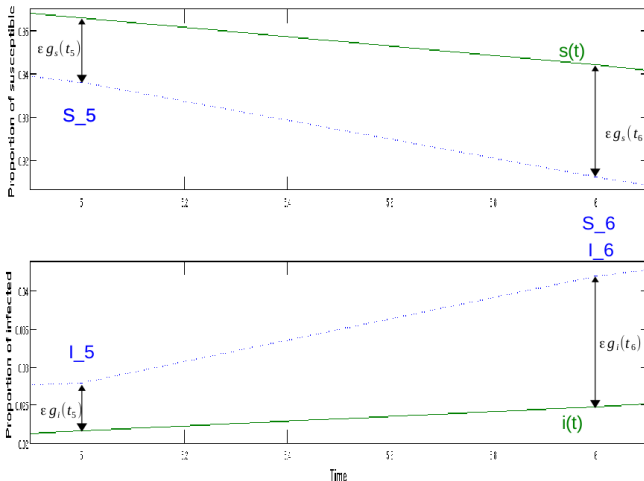


Figure: Diffusion (blue),  $x_{\theta_1}(t)$  (green)

What kind of distance is minimized: comparison with Least squares

Zoom between  $t = 5$  and  $t = 6$

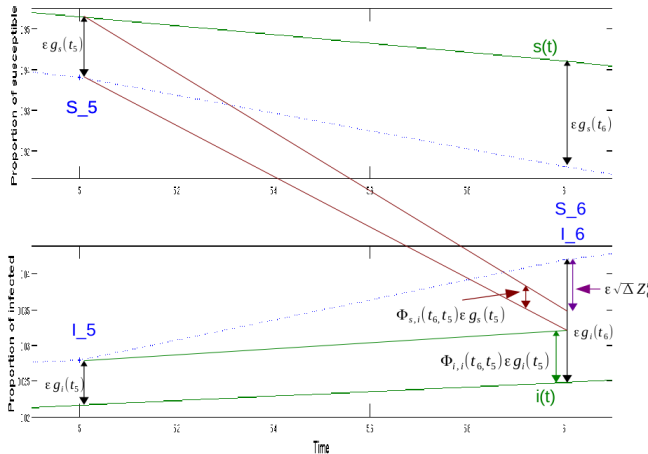
$$g_{\eta}(t_k) = \Phi_{\theta_1}(t_k, t_{k-1})g_{\eta}(t_{k-1}) + \sqrt{\Delta}Z_k^{\eta}$$



## What kind of distance is minimized: comparison with Least squares

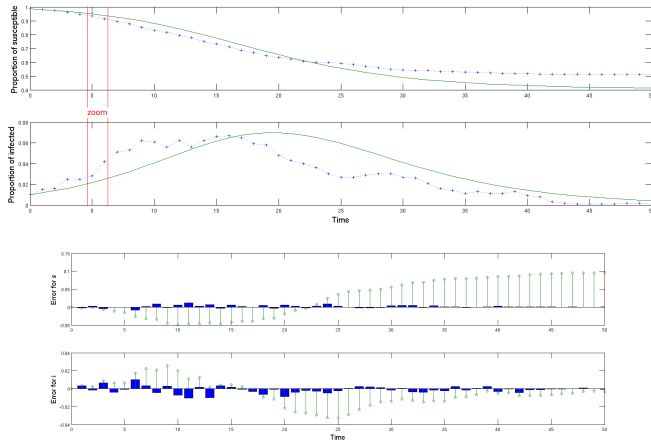
Zoom between  $t = 5$  and  $t = 6$ 

$$g_\eta(t_k) = \Phi_{\theta_1}(t_k, t_{k-1})g_\eta(t_{k-1}) + \sqrt{\Delta}Z_k^\eta$$



## What kind of distance is minimized: comparison with Least square

Figure: Distance to the deterministic model

Figure: Comparison:  $N_k(X, \theta_1)$  (blue) and  $X_{t_k} - x_{\theta_1}(t_k)$  (green)

Results about  $dX_t = b(\theta_1, X_t)dt + \epsilon\sigma(\theta_2, X_t)dB_t$

Under classical regularity assumptions on  $b$  and  $\sigma$ , we proved

$n$  fixed,  $\epsilon \rightarrow 0$ , low frequency contrast

Identifiability assumption:  $\theta_1 \neq \theta'_1 \Rightarrow \{\exists k, 1 \leq k \leq n, x_{\theta_1}(t_k) \neq x_{\theta'_1}(t_k)\}$ .

1. General case (no information on  $\theta_2$ )

$$\epsilon^{-1} (\bar{\theta}_{1\epsilon} - \theta_1^0) \xrightarrow{\epsilon \rightarrow 0} \mathcal{N}(0, J_{\Delta}^{-1}(\theta_1^0, \theta_2^0))$$

2. case  $\theta_2 = f(\theta_1)$  (with information on  $\theta_2$ )

$$\epsilon^{-1} (\tilde{\theta}_{1\epsilon} - \theta_1^0) \xrightarrow{\epsilon \rightarrow 0} \mathcal{N}(0, I_{\Delta}^{-1}(\theta_1^0, \theta_2^0))$$

3.  $n \rightarrow \infty$   $\epsilon \rightarrow 0$ : high frequency contrast

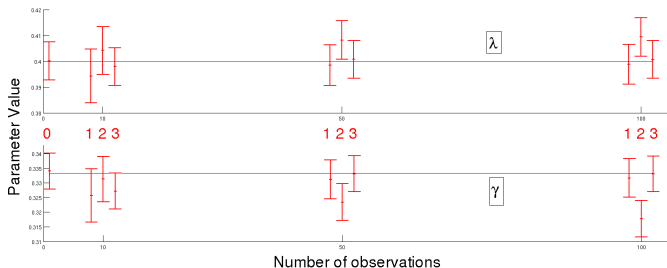
$$\left( \begin{array}{c} \epsilon^{-1} (\check{\theta}_{1\epsilon, \Delta} - \theta_1^0) \\ \sqrt{n} (\check{\theta}_{2\epsilon, \Delta} - \theta_2^0) \end{array} \right) \xrightarrow{n \rightarrow \infty, \epsilon \rightarrow 0} N \left( 0, \left( \begin{array}{cc} I_b^{-1}(\theta_1^0, \theta_2^0) & 0 \\ 0 & I_{\sigma}^{-1}(\theta_1^0, \theta_2^0) \end{array} \right) \right)$$

Remarks

- Epidemics:  $\epsilon = \frac{1}{\sqrt{N}}$ ,  $\theta = \theta_1 = \theta_2$ , then for contrast 3:  $I_b$  is the same as for the Markov jump process (all jumps observed)
- $J_{\Delta}$  is not optimal, but  $I_{\Delta}$  is, in the sense that  $I_{\Delta}(\theta_1, \theta_2) \xrightarrow{\Delta \rightarrow 0} I_b(\theta_1, \theta_2)$

Results on  $SIR$  for  $N = 10000$ , empirical mean estimators on 1000 runs and 95% theoretical CI  
 ( $R_0 = 1.2$ ,  $d = 3$ )

Figure: 0: MLE, 1:  $\bar{\theta}_{1,\epsilon}$  low frequency MCE (general case), 2:  $\tilde{\theta}_{1,\epsilon}$  low frequency MCE ( $\theta_1 = \theta_2$ ), 3:  $\hat{\theta}_{1,\epsilon,\Delta}$  high frequency MCE



### About unrepresented results

- Good results even for  $N = 100$  (our methods seem more robust than MLE)
- Similar performance (w.r.t. MLE) on more sophisticated models

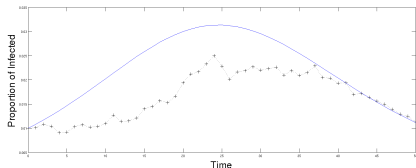
About  $\tilde{\theta}_{1\epsilon}$  (Low frequency with information on  $\theta_2$ ) $\theta_2$  unknown

$$1. \bar{U}_\epsilon(\theta_1) = \frac{1}{\Delta} \sum_{k=1}^n {}^t N_k(\theta_1, X) N_k(\theta_1, X)$$

With information on  $\theta_2$ 

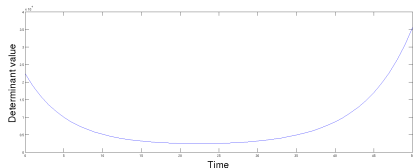
$$2. \tilde{U}_\epsilon(\theta_1) = \frac{1}{\Delta} \sum_{k=1}^n {}^t N_k(\theta_1, X) (\tilde{S}_k^{\theta_1})^{-1} N_k(\theta_1, X)$$

Figure: Comparison between Data and ODE  
 $(x_{\tilde{\theta}_1}(t))$



Not good fit of the data

Figure: Evolution of  $\det(\Sigma^{-1}(\theta_1^0, x_{\theta_1^0}(t)))$



Too much weight on the boundaries



About  $\tilde{\theta}_{1,\epsilon}$  (Low frequency with information on  $\theta_2$ )

Comparing to high frequency contrast

$$3. \check{U}_{\Delta, \epsilon}(\theta_1, \theta_2) = \epsilon^2 \sum_{k=1}^n \log [\det(\Sigma_k)] + \frac{1}{\Delta} \sum_{k=1}^n {}^t N_k(\theta_1) \Sigma_k^{-1} N_k(\theta_1)$$

where  $\Sigma_k = \Sigma(\theta_2, X_{t_{k-1}})$

Corrected contrast with information on  $\theta_2$

$$2'. \check{U}_{\epsilon}^{cor}(\alpha) = \epsilon^2 \sum_{k=1}^n \log [\det(\check{S}_k^{\theta_1})] + \frac{1}{\Delta} \sum_{k=1}^n {}^t N_k(\theta_1) (\check{S}_k^{\theta_1})^{-1} N_k(\theta_1)$$

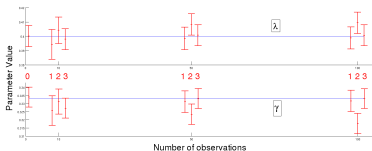


Figure: Previous results

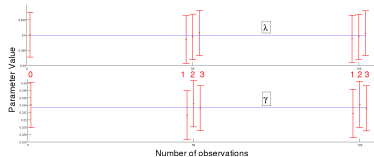
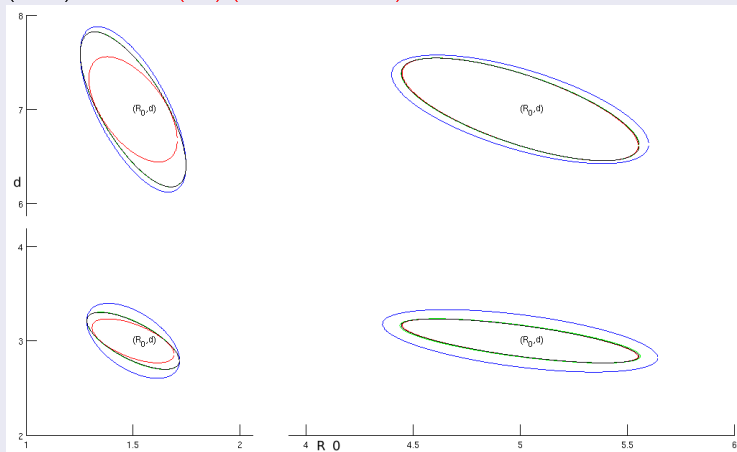


Figure: Corrected results

## Different shapes of the confidence ellipsoids

## SIR Model confidence ellipsoids for the corrected contrast

 $N = 1000, (R_0, d) = \{(1.5, 3), (1.5, 7), (5, 3), (5, 7)\}$ 

 Number of observations:  $n = 10$  (blue),  $n = 1000$  (green),  $n = 2000$  (black), MLE CR (red) (Theoretical limit)


Bias of MLE ( $N = 400$ ;  $R_0 = 1.5$ ;  $d = \{3, 7\}$ ;  $(s_0, i_0) = (0.99, 0.01)$ )

Too much variability

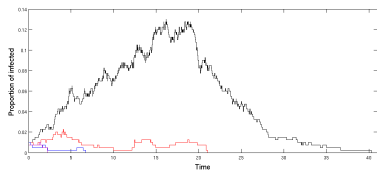
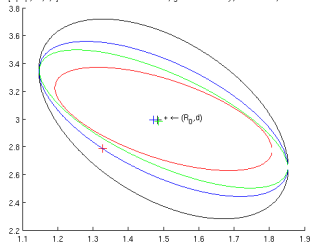


Figure: Trajectories for  $d=3$

Need of an empiric threshold: 5% of infected

Estimation results ( $d = 3$ )

[Npop,R0,d,T]=400 1.5 3 40 : red = limitMLE, green=1obs/day, blue: n=10, black n=5

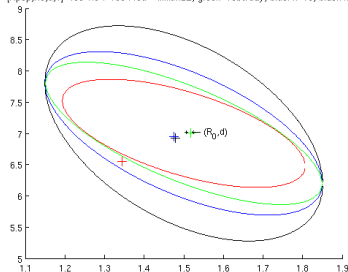


Other values of  $d$  were investigated

Other thresholds: time of extinction, number max of infected

Results ( $d = 7$ , time of extinction)

[Npop,R0,d,T]=400 1.5 7 100 : red = limitMLE, green=1obs/day, blue: n=10, black n=5



Zoom on the red trajectory (see Figure)

$$(R_0, d)_{MLE} = (0.8749, 2.9945)$$

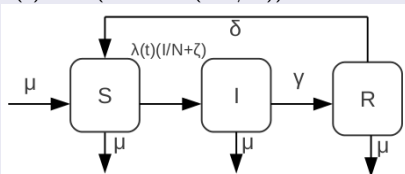
$$(R_0, d)_{LSE} = (0.94, 2.5794)$$

$$(R_0, d)_{cont} = (1.01, 3.0973)$$

Temporal dependence (SIRS) :  $\lambda_1$  difficult to estimate

SIRS with constant immigration in I class

$$\lambda(t) = \lambda_0(1 + \lambda_1 \sin(2\pi t/T))$$



Values

$$R_0 = 1.5; d = 3d; \frac{1}{\delta T_{per}} = 2\gamma,$$

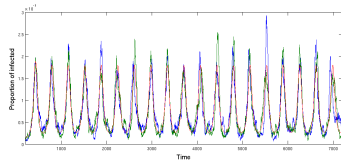
$$\lambda_1 = \{0.05, 0.15\}$$

$$\text{Fixed: } T_{per} = 365, \mu = 1/50 T_{per}, \zeta = \frac{10}{N},$$

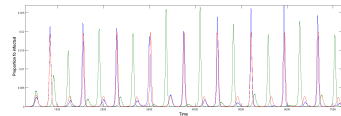
$$N = 10^7, 1 \text{ obs/day(week) for 20 years}$$

- Term for  $\lambda_1$  in  $I_b(\theta_0)$  very small  $\Rightarrow N > 10^5$  for satisfactory CI
- $\lambda_1$  bifurcation parameter for the ODE

- $\lambda_1 = 0.05$  (weak seasonality)



- $\lambda_1 = 0.15$  (stronger seasonality)



Detailed results not shown

Main idea:  $R_0, d, \delta$  well estimated  
 $\lambda_1$ : biased (often estimated to 0)

## Outline

- 1 Characteristics of the epidemic process
  - Constraints imposed by the observation of the epidemic process
  - Simple mechanistic models
- 2 Various mathematical approaches for epidemic spread
  - Natural approach: Markov jump process
  - First approximation by ODEs
  - Gaussian approximation of the Markov jump process
  - Diffusion approximation of the Markov jump process
- 3 Inference for discrete observations of diffusion or Gaussian processes with small diffusion coefficient
  - Contrast processes for fixed or large number of observations
  - Correction of a non asymptotic bias
  - Comparison of estimators on simulated epidemics
- 4 Epidemics incompletely observed: partially and integrated diffusion processes (Work in progress)
  - Back to epidemic data
  - Inference approach: Work in progress

Incidence for SIR models

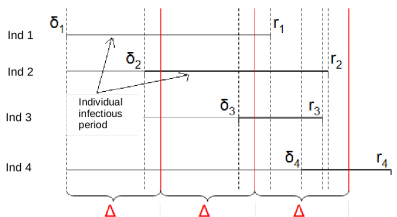
Discrete observation of all the coordinates

- Confined studies
- Childhood diseases



SIR Model:

- Incidence at  $t_2$ :  $\int_{t_1}^{t_2} \lambda S(t) \frac{I(t)}{N} dt$
- High frequency data  $\approx \lambda S(t_2) \frac{I(t_2)}{N}$
- $d$  small : new infected  $\approx$  new removed,  
 $\int_{t_1}^{t_2} \gamma I(t) dt = R(t_2) - R(t_1)$



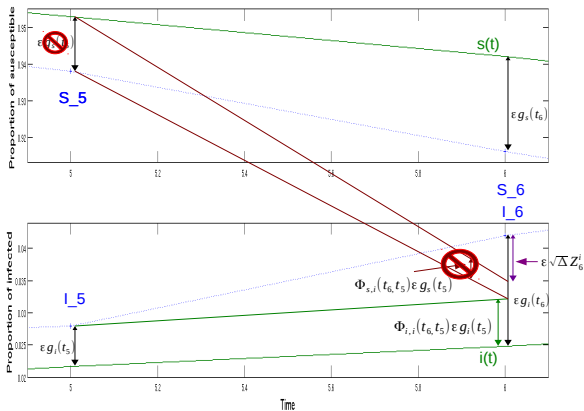
Individual Tracking	{ Ind 1, Ind 2 }	{Ind 1, Ind 2, Ind 3 }	{Ind 4 }
Number of Infected	2	3	1
Incidence (new infected)	2	1	1

**Diffusion perspectives**

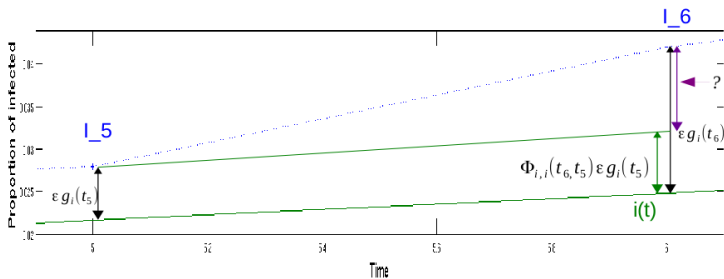
- Partial and discrete obs. of the diffusion process (Itô's formula)
- Partial and Integrated discrete obs.

## Partially observed diffusion process: initial idea

Previous main idea not directly applicable:



## Partially observed diffusion process: initial idea

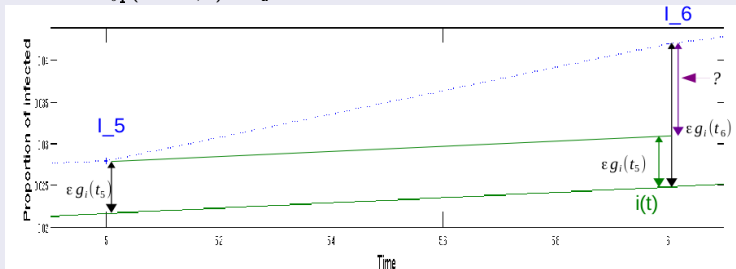


No good properties for  $n$  fixed,  $\epsilon \rightarrow 0$ .



## Partially observed diffusion process: initial idea

Use that  $\Phi_{\theta_1}(t + \Delta, t) \approx I_d$  as  $\Delta \rightarrow 0$



Function of the observations :  $I_{t_k} - i(t_k) - I_{t_{k-1}} + i(t_{k-1})$

## Integrated diffusion process

Integration of the relation :  $g_{\eta}(t_k) = \Phi_{\theta_1}(t_k, t_{k-1})g_{\eta}(t_{k-1}) + \sqrt{\Delta}Z_k^{\eta}$

$\Rightarrow$  link between  $g$  and the integrated process : similar to Kalman filtering techniques