

# Monte Carlo adjusted profile likelihood, with applications to spatiotemporal and phylodynamic inference.

Edward Ionides

University of Michigan, Department of Statistics

Isaac Newton Institute workshop on

*Future challenges in statistical scalability*

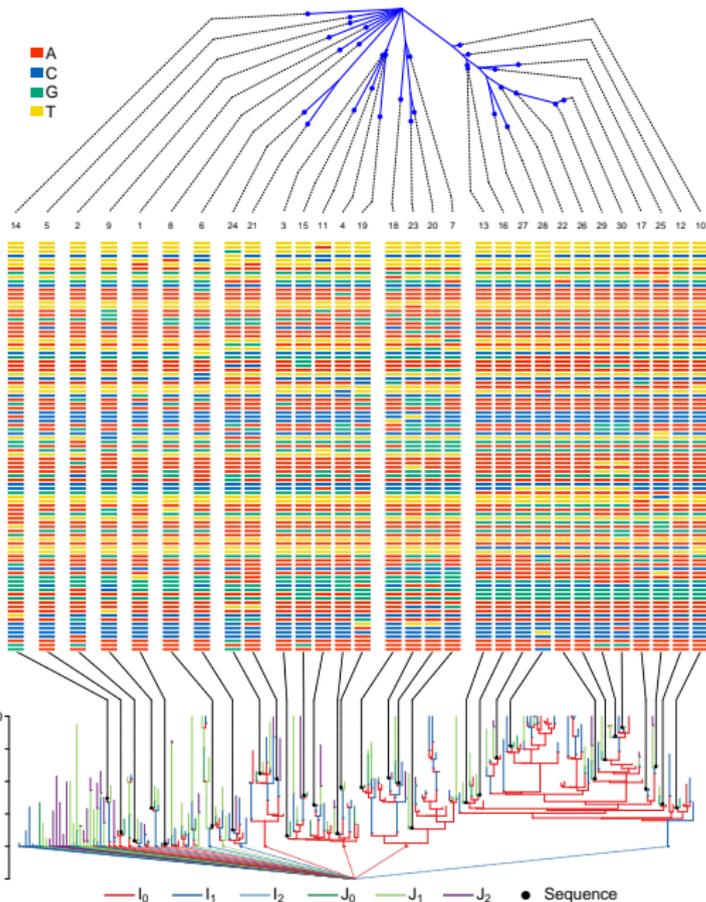
Thursday 28th June, 2018

Collaborators: Carles Bretó, Aaron King, Joonha Park, Alex Smith

Slides are online at

<http://dept.stat.lsa.umich.edu/~ionides/talks/ini18.pdf>

# HIV transmission dynamics: inference for a complex system with not-small data (Smith et al., 2017)



Top: phylogeny of “observed”  
sequences from a simulation.

Middle: simulated sequences.  
Actual data were 100 partial  
HIV sequences of length  $\approx$   
1000.

Bottom: Transmission forest  
for the full epidemic.

red: undiagnosed early infection  
blue: undiagnosed chronic infection  
green: diagnosed

Our agenda is to collect together some approaches to likelihood based inference:

- Profile likelihood.
- Fisher information (observed and expected).
- Local asymptotic normality.
- Smoothed likelihood.
- Monte Carlo likelihood.

We then see how these ideas combine to facilitate inference for some complex dynamic systems.

## Profile likelihood: some definitions

- The log likelihood function for model  $f_Y(y; \theta)$  and data  $y^*$  is  $\lambda(\theta; y^*) = \log f_Y(y^*; \theta)$ ,
- A maximum likelihood estimate (MLE) is  $\hat{\theta}^* = \hat{\theta}(y^*) = \arg \max_{\theta} \lambda(\theta; y^*)$ .
- We suppose  $\theta = (\phi, \psi)$  with  $\phi \in \mathbb{R}^1$  and  $\psi \in \mathbb{R}^{p-1}$ . Here,  $\phi$  is a focal parameter for which we are interested in obtaining a confidence interval.
- The **profile log likelihood function** for  $\phi$  is defined as  $\lambda^P(\phi; y^*) = \max_{\psi} \lambda((\phi, \psi); y^*)$ .
- The profile log likelihood is maximized at a marginal MLE,  $\hat{\phi}^* = \hat{\phi}(y^*) = \arg \max_{\phi} \lambda^P(\phi; y^*)$ .
- A profile likelihood confidence interval with cutoff  $\delta$  is defined as  $\{\phi : \lambda^P(\phi; y^*) \geq \lambda^P(\hat{\phi}^*; y^*) - \delta\}$ .

## Profile likelihood: some history

- Profile likelihood confidence intervals are equivalent to likelihood ratio tests, which have a long history.
- Box and Cox (1964) graphed the profile likelihood under the name of maximized likelihood, and constructed confidence intervals using the  $\chi^2$  cutoff.
- Cox and Snell (1970) made an early use of the name “profile likelihood.” Use of “maximized likelihood” continued through 1970’s but is now antiquated.
- Much work in the 1980’s focused on how to modify profile likelihood for improved higher-order asymptotic behavior (Barndorff-Nielsen, 1983).
- Profile likelihood has uses in semiparametric inference (Murphy and van der Vaart, 2000). The proportional hazard “partial likelihood” (Cox, 1972) is a semiparametric profile likelihood.

# Fisher information and observed Fisher information

- Fisher information, evaluated at the MLE, is

$$I_{ij} = \mathbb{E} \left[ -\frac{\partial}{\partial \theta_i \partial \theta_j} \lambda(\hat{\theta}^*; Y) \right]$$

- Observed Fisher information is

$$I_{ij}^* = -\frac{\partial}{\partial \theta_i \partial \theta_j} \lambda(\hat{\theta}^*; y^*)$$

- The asterisk denoting observed Fisher information indicates the additional data dependence.
- Corresponding standard errors and 95% confidence intervals for  $\phi = \theta_i$  are

$$\text{SE}_F = \sqrt{[I^{-1}]_{ii}} \quad \text{CI}_F = [\hat{\theta}^* - 1.96 \text{SE}_F, \hat{\theta}^* + 1.96 \text{SE}_F]$$

$$\text{SE}_F^* = \sqrt{[I^{*-1}]_{ii}} \quad \text{CI}_F^* = [\hat{\theta}^* - 1.96 \text{SE}_F^*, \hat{\theta}^* + 1.96 \text{SE}_F^*]$$

- An identity:  $\frac{d^2}{d\phi^2} \lambda_P(\phi; y^*) = -\left[ [I^{*-1}]_{ii} \right]^{-1}$ .
- For a quadratic likelihood function,  $\text{CI}_F^*$  is equal to the profile likelihood confidence interval,  
 $\text{CI}_P = \{ \phi : \lambda_P(\phi; y^*) \geq \lambda_P(\hat{\phi}^*; y^*) - 1.92 \}$ .

# In favor of observed Fisher information and profile likelihood

- Heuristically, the error on an estimator depends on the amount of information observed in the actual experiment.
- Efron and Hinkley (1978) argued for the observed Fisher standard error,  $SE_F^*$ , over  $SE_F$ . Formal reasoning was limited to special cases, with arguments based on ancillarity.
- Lindsay and Li (1997) used a risk framework to show  $SE_F^{*2}$  gives an asymptotic 2nd order mean square optimal estimate of  $(\hat{\phi} - \phi)^2$ , unlike  $SE_F^2$  or the bootstrap.
- $CI_P$  transforms naturally if  $\phi$  is reparameterized by a monotonic  $h(\phi)$ . If there is an unknown  $h$  such that the log likelihood is [approximately] quadratic for any  $y^*$ ,  $CI_P$  [approximately] corresponds to  $CI_F^*$  computed on this scale. Thus, heuristically, we may expect  $CI_P$  to have comparable asymptotic optimality to  $CI_F^*$  but better finite sample behavior.

## Local asymptotic normality

- LAN (Le Cam, 1986) concerns a sequence of statistical models,  $f_{Y,n}(y_n; \theta)$ , and the behavior of the log likelihood ratio,  $\Lambda_n(\theta) = \log f_{Y_n}(Y_n; \theta) - \log f_{Y_n}(Y_n; \theta_0)$  when  $Y_n \sim F_{Y,n}(y_n; \theta_0)$ .
- $f_{Y,n}(y_n; \theta)$  has LAN with information matrix  $K$  if there is a sequence of random variables  $\Delta_n \xrightarrow{d} N(0, K)$  such that, for all bounded  $\{t_n\}$ ,

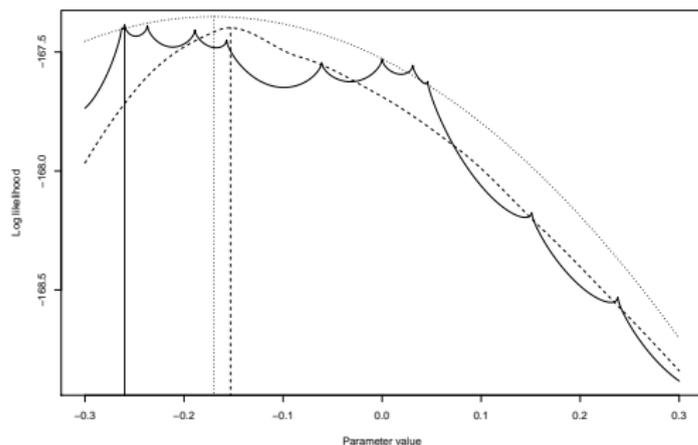
$$\Lambda_n(\theta_0 + t_n n^{-1/2}) = t_n^T \Delta_n - \frac{1}{2} t_n^T K t_n + o_p(1; \theta_0).$$

- $K$  is the asymptotic information rate concerning  $\theta$ ; it coincides with the Fisher information under regularity conditions. From Hájek's convolution theorem, an estimator  $\hat{\theta}_n$  is asymptotically efficient if LAN holds and  $n^{1/2}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N(0, K^{-1})$
- Bickel et al. (1993) discuss LAN and demonstrate the utility of the LAN framework for semiparametrics

# Maximum quadratic likelihood estimator (MQLE)

- LAN justifies the following estimation procedure:
  - ① Evaluate the log likelihood at a grid of points in the neighborhood of a  $\sqrt{n}$ -consistent estimator.
  - ② Fit a quadratic through these points.
  - ③ Obtain the maximum of this quadratic.
- This is called Le Cam's one-step estimator. We will say **maximum quadratic likelihood estimator (MQLE)**.
- MQLE is efficient under LAN and more generally when the likelihood is locally asymptotic quadratic (Le Cam, 1986).
- Under regularity, LAN is equivalent to asymptotic normality of MLE.
- MQLE can succeed when LAN holds but the MLE behaves badly.
- LAN can be easier to prove than asymptotic normality of the MLE.
- The one-step estimator is, in some sense, better than the MLE.
- **But, who would use LAN for data analysis if the log likelihood is grossly non-quadratic**

A non-smooth likelihood: An iid model,  $Y = Y_1, \dots, Y_n$  with  $f_Y(y|\theta) \propto \exp \left\{ - \sum_{i=1}^n |y_i - \theta|^\omega \right\}$



- Dotted: MQLE, initialized at the median.
- Dashed: maximum smoothed likelihood estimator (MSLE).
- Solid: likelihood and MLE.
- Truth:  $\theta = 0$ ,  $\omega = 0.6$  and  $n = 100$ .

- This model does not satisfy the usual Cramér conditions for the MLE. MQLE and MSLE are 15% more efficient than MLE (Ionides, 2005).
- Perhaps more importantly, they are not worse!
- **For difficult likelihood surfaces, or when we must rely on Monte Carlo approximation of the likelihood, MQLE and MSLE may be easier to implement than MLE.**

# More on maximum smoothed likelihood estimation (MSLE)

- MSLE (Ionides, 2005) involves the following steps:
  - ① Evaluate the log likelihood at a grid of points in the neighborhood of a  $\sqrt{n}$ -consistent estimator.
  - ② Fit a smooth curve through these points.
  - ③ Obtain the maximum of this smooth curve.
- MSLE replaces the quadratic of MQLE with a smoother.
- The smoothed likelihood can be used to construct profile confidence intervals.
- As long as the smoother fits a quadratic through points on a quadratic, MSLE inherits asymptotic optimality from MQLE.
- The loess smoother in R is a 2nd order local polynomial smoother with this property.

# Monte Carlo profile confidence intervals for dynamic systems

- Monte Carlo methods to evaluate and maximize the likelihood function enable the construction of confidence intervals and hypothesis tests, facilitating scientific investigation using models for which the likelihood function is intractable.
- When Monte Carlo error can be made small, by sufficiently exhaustive computation, then the standard theory and practice of likelihood-based inference applies. One may still want to use MSLE to enable reliable inference at reduced computational cost.
- As datasets become larger, and models more complex, situations arise where no reasonable amount of computation can render Monte Carlo error negligible.
- We seek profile likelihood methodology enabling frequentist inferences accounting for Monte Carlo error.
- This methodology facilitates inference for computationally challenging dynamic latent variable models.

# A metamodel for a Monte Carlo profile

- A **Monte Carlo metamodel** is a statistical model fitted to output of a Monte Carlo algorithm.
- We have independent Monte Carlo profile likelihood evaluations  $(\check{\lambda}_k^P(y^*), k \in 1:K)$  at points  $\phi_{1:K} = (\phi_1, \dots, \phi_K)$ .
- Without loss of generality we can write

$$[M1] \quad \check{\lambda}_k^P(y^*) = \lambda^P(\phi_k; y^*) + \beta_k(y^*) + \epsilon_k(y^*), \quad k \in 1:K,$$

where Monte Carlo errors  $\epsilon_{1:K}(Y)$  are, **by construction, mean zero and independent** conditional on  $Y$ . In [M1],  $\beta_k(y^*)$  is Monte Carlo bias.

- Local to the MLE, we may make additional metamodel assumptions:  
[M2]  $\beta_k(y^*) = \beta(y^*)$  : constant bias.  
[M3]  $\text{Var}[\epsilon_k(y^*)] = \sigma^2(y^*) < \infty$  : constant variance.
- We can complete the metamodel by proposing parametric or nonparametric specifications of  $\lambda^P(\phi; y^*)$ .

# A quadratic metamodel for the profile likelihood

- LAN suggests a quadratic metamodel,

$$\check{\lambda}_k^P(y) = -\hat{a}(y)\phi_k^2 + \hat{b}(y)\phi_k + \hat{c}(y) + \epsilon_k, \quad \text{Var}(\epsilon_k) = \sigma^2(y).$$

- The unknown coefficients  $\hat{a}^* = \hat{a}(y^*)$ ,  $\hat{b}^* = \hat{b}(y^*)$  and  $\hat{c}^* = \hat{c}(y^*)$  make a quadratic approximation to the intractable likelihood.
- We fit the metamodel to the Monte Carlo profile evaluations, using linear regression to estimate  $(\hat{a}^*, \hat{b}^*, \hat{c}^*)$  by  $(\check{a}^*, \check{b}^*, \check{c}^*)$ .
- The marginal MLE  $\hat{\phi}^*$  can be approximated by the maximum of  $\check{\lambda}^Q(\phi; y^*)$ , which is given by  $\check{\phi}^Q(y^*, \epsilon) = \check{b}(y^*, \epsilon)/2\check{a}(y^*, \epsilon)$
- We can separate the variability of  $\check{\phi}^Q(Y, \epsilon)$  into two components:
  - 1 **Statistical error** is uncertainty from randomness in the data, viewed as a draw from the statistical model. This is the usual statistical error of  $\hat{b}(y^*)/2\hat{a}(y^*)$  as an estimate of  $\phi$ .
  - 2 **Monte Carlo** error is the uncertainty from implementing a Monte Carlo estimator. This is the error in  $\check{b}(y^*, \epsilon)/2\check{a}(y^*, \epsilon)$  as a Monte Carlo estimate of  $\hat{b}(y^*)/2\hat{a}(y^*)$ .

## Monte Carlo error and statistical error

- Routine application of the delta method gives a central limit approximation for the Monte Carlo error on the maximum, conditional on  $Y = y^*$ ,

$$\frac{\check{b}^*}{2\check{a}^*} \approx N \left[ \left( \frac{\hat{b}^*}{2\hat{a}^*} \right), \text{SE}_{\text{mc}}^2 \right],$$

where

$$\text{SE}_{\text{mc}}^2 = \frac{1}{4\check{a}^{*2}} \left\{ \text{Var}[\check{b}^*] - \frac{2\check{b}^*}{\check{a}^*} \text{Cov}[\check{a}^*, \check{b}^*] + \frac{\check{b}^{*2}}{\check{a}^{*2}} \text{Var}[\check{a}^*] \right\}.$$

- The usual statistical standard error,  $1/\sqrt{2\hat{a}^*}$ , is not available to us. Its Monte Carlo estimate is

$$\text{SE}_{\text{stat}} = \frac{1}{\sqrt{2\check{a}^*}}.$$

- Under suitable regularity, these two error sources are additive, and so

$$\text{SE}_{\text{total}} = \sqrt{\text{SE}_{\text{mc}}^2 + \text{SE}_{\text{stat}}^2}.$$

## Using $SE_{\text{total}}$ for a Monte Carlo adjusted profile (MCAP)

- The usual  $\chi^2$  cutoff for profile confidence intervals is based on quadratic asymptotics for the log likelihood. It is robust to reparameterization, and can be applied to either the actual profile or a smoothed version.
- Exactly the same argument can be applied to give a cutoff for a smoothed Monte Carlo profile based on a quadratic approximation:

$$\delta = \check{a}^* \times (z_\alpha \times SE_{\text{total}})^2 = z_\alpha^2 \left( \check{a}^* \times SE_{\text{mc}}^2 + \frac{1}{2} \right),$$

where  $z_\alpha$  is the  $1 - \alpha/2$  normal quantile.

- if  $SE_{\text{mc}} = 0$ , the cutoff for  $\alpha = 0.05$  reduces to  $\delta = 1.96^2/2 = 1.92$ .
- We apply this cutoff after estimating the profile via a locally weighted quadratic smoother.  $SE_{\text{mc}}$  can be computed using the local weights at the maximum.

## Using $SE_{\text{total}}$ for a Monte Carlo adjusted profile (MCAP)

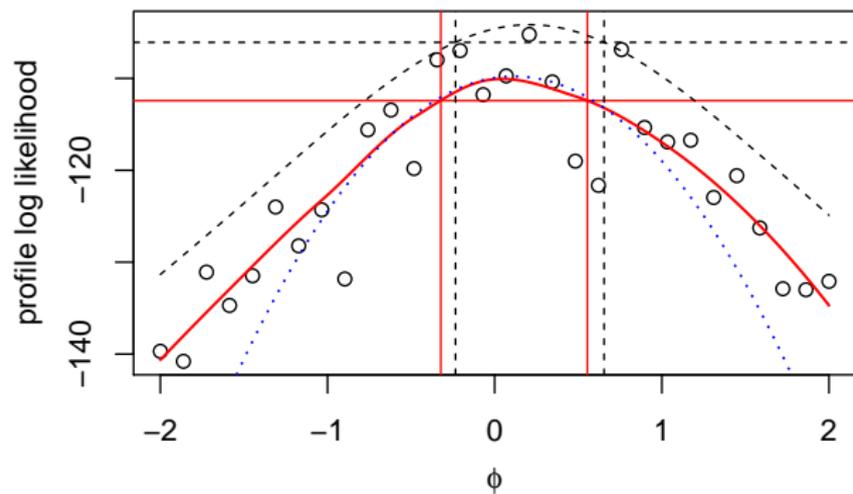
- The usual  $\chi^2$  cutoff for profile confidence intervals is based on quadratic asymptotics for the log likelihood. It is robust to reparameterization, and can be applied to either the actual profile or a smoothed version.
- Exactly the same argument can be applied to give a cutoff for a smoothed Monte Carlo profile based on a quadratic approximation:

$$\delta = \check{a}^* \times (z_\alpha \times SE_{\text{total}})^2 = z_\alpha^2 \left( \check{a}^* \times SE_{\text{mc}}^2 + \frac{1}{2} \right),$$

where  $z_\alpha$  is the  $1 - \alpha/2$  normal quantile.

- if  $SE_{\text{mc}} = 0$ , the cutoff for  $\alpha = 0.05$  reduces to  $\delta = 1.96^2/2 = 1.92$ .
- We apply this cutoff after estimating the profile via a locally weighted quadratic smoother.  $SE_{\text{mc}}$  can be computed using the local weights at the maximum.
- We call this procedure a **Monte Carlo adjusted profile (MCAP)**.

# A toy: importance sampling for a log normal model



Points show Monte Carlo profile evaluations. Black dashed lines: exact profile and 95% confidence interval. Solid red lines: MCAP confidence interval. Dotted blue line: quadratic approximation.

	Exact profile	MCAP profile	Bootstrap	Quadratic
Coverage %	94.3	93.4	93.3	93.3
Mean width	0.78	0.88	0.94	0.92

# Statistical challenges for nonlinear mechanistic modeling in ecology and epidemiology

- 1 Combining measurement noise and process noise.
- 2 Including covariates in mechanistically plausible ways.
- 3 Continuous time models.
- 4 Modeling and estimating interactions in coupled systems.
- 5 Dealing with unobserved variables.
- 6 **Spatiotemporal data and models.**
- 7 **Inferences from genetic sequence data.**

(1–6) were enumerated by Bjornstad and Grenfell (*Science*, 2001).

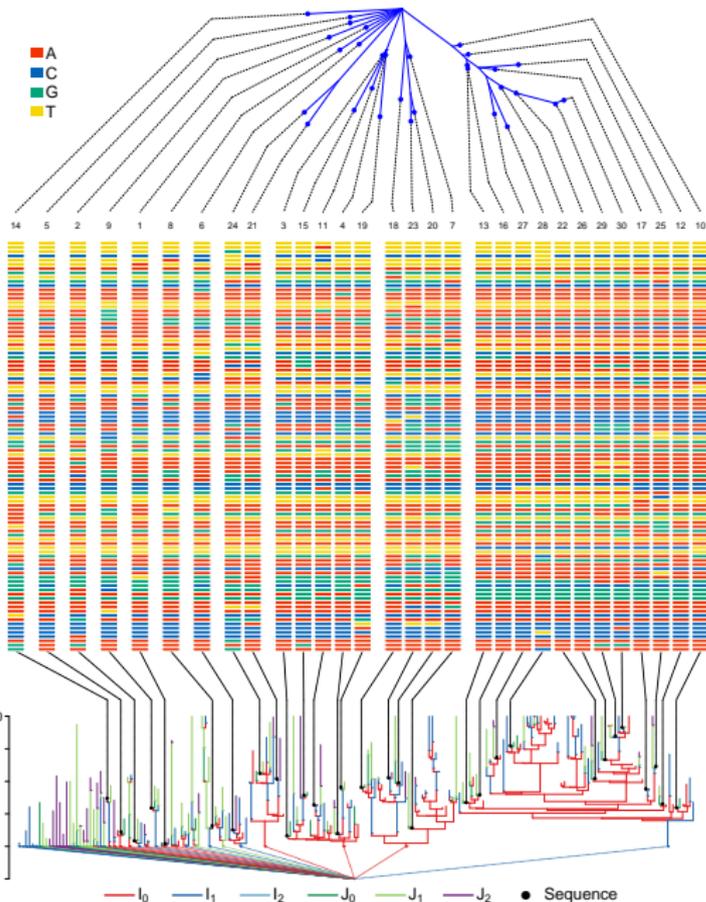
(1–5) are now routinely solved using modern methods for nonlinear partially observed Markov process (POMP) models (Ionides et al., 2015; King et al., 2016).

(7) was described by Grenfell et al (*Science*, 2004) and a general POMP solution was shown by Smith et al (*Molecular Biology & Evolution*, 2017).

# Inferring population dynamics from genetic sequence data

- Genetic sequence data on a sample of individuals in an ecological system has potential to reveal population dynamics.
- Extraction information on population dynamics from genetic data has been termed *phylogenetics* (Grenfell et al., 2004).
- Inference via the full likelihood stretches modern computational capabilities, but can be done using the `genPomp` algorithm of Smith et al. (2017).
- The `genPomp` algorithm is an application of iterated filtering methodology (Ionides et al., 2015) to phylodynamic models and data.
- However, the `genPomp` algorithm leads to estimators with high Monte Carlo variance, indeed, too high for reasonable amounts of computation resources to reduce Monte Carlo variability to negligibility.
- This situation provides a useful scenario to demonstrate our methodology.

# HIV transmission dynamics: inference for a complex system with not-small data (Smith et al., 2017)



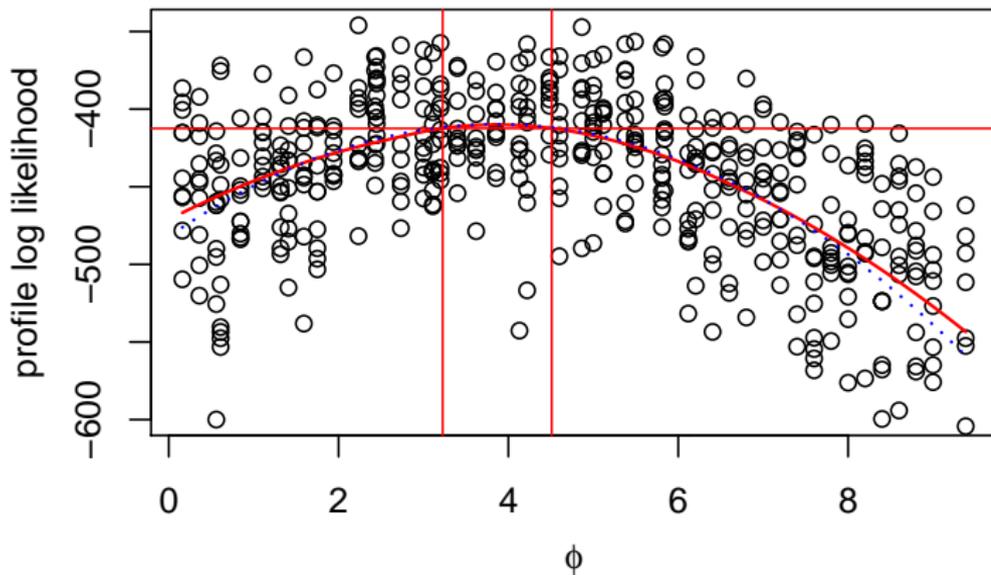
Top: phylogeny of “observed”  
sequences from a simulation.

Middle: simulated sequences.  
Actual data were 100 partial  
HIV sequences of length  $\approx$   
1000.

Bottom: Transmission forest  
for the full epidemic.

red: undiagnosed early infection  
blue: undiagnosed chronic infection  
green: diagnosed

# Monte Carlo profile for genetic data on HIV dynamics

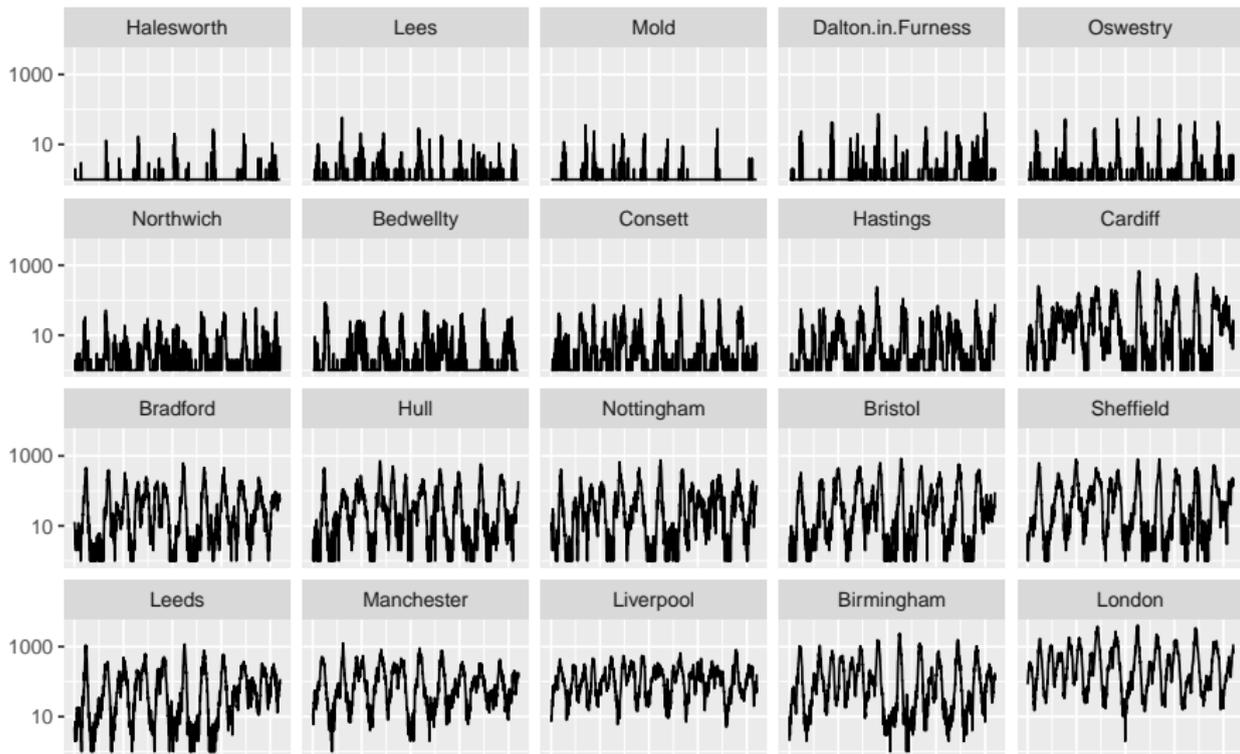


- $\phi$  models HIV transmitted by recently infected, diagnosed individuals.
- The MCAP cutoff is 2.35, compared to the unadjusted cutoff of 1.92.
- The computation of this figure took approximately 10 days using 500 cores on a Linux cluster.
- The standard error of the profile evaluations is around 25 log units.

# Inference for nonlinear partially observed spatiotemporal systems

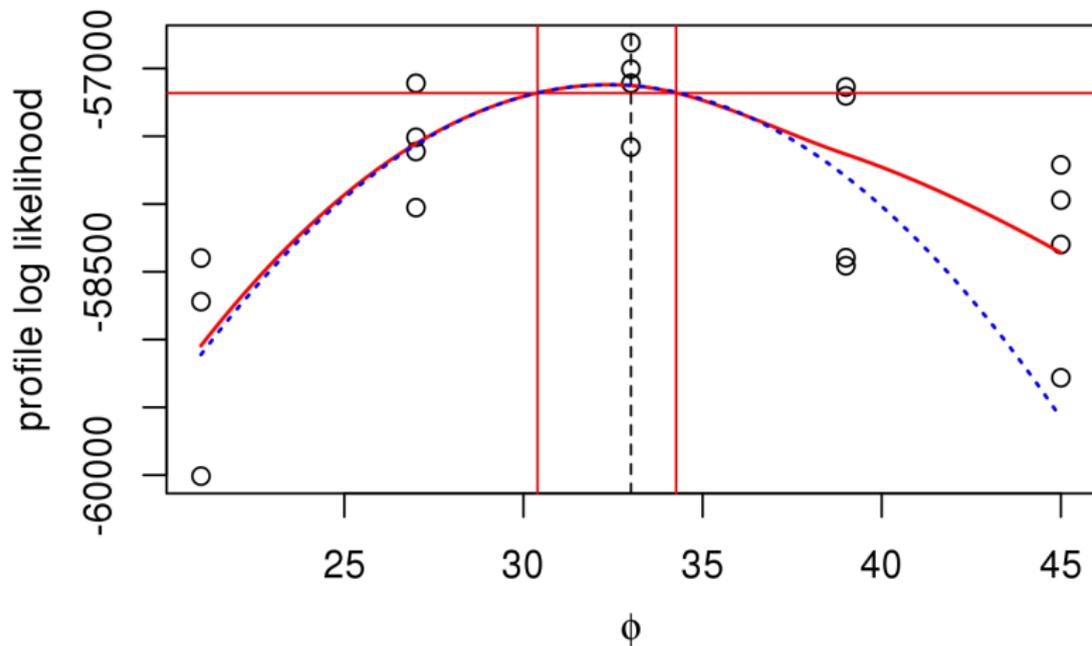
- Sequential Monte Carlo (SMC) methods that enable routine data analysis in low-dimensional systems scale poorly for higher dimensions (Bengtsson et al., 2008).
- Under weak coupling assumptions, localized SMC can theoretically succeed in higher dimensions (Rebeschini and van Handel, 2015).
- We used an SMC algorithm, the Guided Intermediate Resampling Filter (GIRF) of Park and Ionides (2017), to study spatiotemporal infectious disease transmission, fitting models to measles case reports.
- The SMC algorithm was coerced into likelihood maximization using the iterated filtering algorithm of Ionides et al. (2015).

# Measles in 20 UK cities, 1944–1965



- Modeled using coupled over-dispersed Markov chains representing susceptible, latent, infectious and recovered individuals.

## Coupled measles SEIR in 20 cities: profiling contact rate



Monte Carlo adjusted profile (MCAP) methodology gives a cutoff of 61.6, rather than the usual 1.92, for the confidence interval construction. Here, Monte Carlo variability is larger than statistical uncertainty. This is a simulation study, with the truth at the vertical dashed line.

## Comparison with methods based on summary statistics

- We have focused on likelihood-based confidence intervals.
- An alternative to likelihood-based inference is to compare the data with simulations using some summary statistic.
- Various plug-and-play methodologies of this kind have been proposed, such as synthetic likelihood (Wood, 2010) and nonlinear forecasting (Ellner et al., 1998).
- For large nonlinear systems, it can be hard to find low-dimensional summary statistics that capture a good fraction of the information in the data.
- Even summary statistics derived by careful scientific or statistical reasoning have been found surprisingly uninformative compared to the whole data likelihood in both scientific investigations (Shrestha et al., 2011) and simulation experiments (Fasiolo et al., 2016).

## Comparison with Bayesian computation

- Much attention has been given to scaling Bayesian computation to complex models and large data. Latent process models are closely related computationally to Bayesian inference: Bayesian parameters are latent random variables.
- Bayesian Numerical methods such as expectation propagation (EP), variational Bayes, and posterior interval estimation (PIE) are effective for some model classes. They emphasize hierarchical models, where the joint density of the data and latent variables can be conveniently factorized. The `genPomp` and spatiotemporal examples don't have this structure: the MCAP methodology has no such requirement.
- Some simulation-based Bayesian methods use unbiased Monte Carlo likelihood evaluations inside an MCMC algorithm (Andrieu and Roberts, 2009). Error in likelihood evaluation slows MCMC convergence. Optimal trade-off between number of MCMC iterations and time spent on each likelihood evaluation occurs at a Monte Carlo likelihood std. deviation of one log unit (Doucet et al., 2015). For our examples, Monte Carlo errors that small are infeasible.

# Conclusions

- MCAP provides a simple and general approach to inference when the signal-to-noise ratio in the Monte Carlo profile log likelihood is sufficient to uncover its main features, up to an unimportant vertical shift.
- For large datasets in which the signal (quantified as the curvature of the log likelihood) is sufficient, the methodology can be effective even when the Monte Carlo noise is far too big to carry out standard Bayesian MCMC techniques.
- Various extensions to the theory and practice of Monte Carlo adjusted likelihood-based inference would be useful for future applied work on large and complex systems.

# The IF2 algorithm (Ionides et al., 2015). Input and output.

## input:

Simulator for latent process initial density,  $f_{X_0}(x_0; \theta)$

**Simulator for transition density**,  $f_{X_n|X_{n-1}}(x_n | x_{n-1}; \theta)$ ,  $n$  in  $1:N$

Evaluator for measurement density,  $f_{Y_n|X_n}(y_n | x_n; \theta)$ ,  $n$  in  $1:N$

Data,  $y_{1:N}^*$

Number of iterations,  $M$

Number of particles,  $J$

Initial parameter swarm,  $\{\Theta_j^0, j \text{ in } 1:J\}$

Perturbation density,  $h_n(\theta | \varphi; \sigma)$ ,  $n$  in  $1:N$

Perturbation sequence,  $\sigma_{1:M}$

**output:** Final parameter swarm,  $\{\Theta_j^M, j \text{ in } 1:J\}$

---

Algorithms that specify the dynamic model via a simulator are said to be **plug-and-play**. This property ensures applicability to the broad class of models for which a simulator is available.

## IF2: iterated SMC with perturbed parameters

For  $m$  in  $1:M$  [ $M$  filtering iterations, with decreasing  $\sigma_m$ ]

$$\Theta_{0,j}^{F,m} \sim h_0(\cdot | \Theta_j^{m-1}; \sigma_m) \text{ for } j \text{ in } 1:J$$

$$X_{0,j}^{F,m} \sim f_{X_0}(x_0; \Theta_{0,j}^{F,m}) \text{ for } j \text{ in } 1:J$$

For  $n$  in  $1:N$  [SMC with  $J$  particles]

$$\Theta_{n,j}^{P,m} \sim h_n(\cdot | \Theta_{n-1,j}^{F,m}, \sigma_m) \text{ for } j \text{ in } 1:J$$

$$X_{n,j}^{P,m} \sim f_{X_n|X_{n-1}}(x_n | X_{n-1,j}^{F,m}; \Theta_j^{P,m}) \text{ for } j \text{ in } 1:J$$

$$w_{n,j}^m = f_{Y_n|X_n}(y_n^* | X_{n,j}^{P,m}; \Theta_{n,j}^{P,m}) \text{ for } j \text{ in } 1:J$$

Draw  $k_{1:J}$  with  $\mathbb{P}(k_j = i) = w_{n,i}^m / \sum_{u=1}^J w_{n,u}^m$

$$\Theta_{n,j}^{F,m} = \Theta_{n,k_j}^{P,m} \text{ and } X_{n,j}^{F,m} = X_{n,k_j}^{P,m} \text{ for } j \text{ in } 1:J$$

End For

$$\text{Set } \Theta_j^m = \Theta_{N,j}^{F,m} \text{ for } j \text{ in } 1:J$$

End For

## IF2 as an iterated Bayes map

- Each iteration of IF2 is a Monte Carlo approximation to a map

$$T_\sigma f(\theta_N) = \frac{\int \check{\ell}(\theta_{0:N}) h(\theta_{0:N} | \varphi; \sigma) f(\varphi) d\varphi d\theta_{0:N-1}}{\int \check{\ell}(\theta_{0:N}) h(\theta_{0:N} | \varphi; \sigma) f(\varphi) d\varphi d\theta_{0:N}}, \quad (1)$$

where  $\check{\ell}(\theta_{0:N})$  is the likelihood of the data under the extended model with time-varying parameter  $\theta_{0:N}$ .

- $f$  and  $T_\sigma f$  in (1) approximate the initial and final density of the IF2 parameter swarm.
- When the standard deviation of the parameter perturbations is held fixed at  $\sigma_m = \sigma > 0$ , IF2 is a Monte Carlo approximation to  $T_\sigma^M f(\theta)$ .
- Iterated Bayes maps are not usually contractions.
- We study the homogeneous case,  $\sigma_m = \sigma$ .
- Studying the limit  $\sigma \rightarrow 0$  may be as appropriate as an asymptotic analysis to study the practical properties of a procedure such as IF2, with  $\sigma_m$  decreasing down to some positive level  $\sigma > 0$  but never completing the asymptotic limit  $\sigma_m \rightarrow 0$ .

**Theorem 1.** Assuming adequate regularity conditions, there is a unique probability density  $f_\sigma$  with

$$\lim_{M \rightarrow \infty} T_\sigma^M f = f_\sigma,$$

with the limit taken in the  $L^1$  norm. The SMC approximation to  $T_\sigma^M f$  converges to  $T_\sigma^M f$  as  $J \rightarrow \infty$ , uniformly in  $M$ .

- Theorem 1 follows from existing results on filter stability.
- Convergence and stability of the ideal filter (a small error at time  $t$  has diminishing effects at later times) is closely related to convergence of SMC.

**Theorem 2.** Under regularity conditions,  $\lim_{\sigma \rightarrow 0} f_\sigma$  approaches a point mass at the maximum likelihood estimate (MLE).

### Outline of proof.

- Trajectories in parameter space which stray away from the MLE are down-weighted by the Bayes map relative to trajectories staying close to the MLE.
- As  $\sigma$  decreases, excursions any fixed distance away from the MLE require an increasing number of iterations and therefore receive an increasing penalty from the iterated Bayes map.
- Bounding this penalty proves the theorem.

**Thank you!**

Slides are online at

<http://dept.stat.lsa.umich.edu/~ionides/talks/ini18.pdf>

## References I

- Andrieu, C. and Roberts, G. O. (2009). The pseudo-marginal approach for efficient computation. *Annals of Statistics*, 37:697–725.
- Barndorff-Nielsen, O. (1983). On a formula for the distribution of the maximum likelihood estimator. *Biometrika*, 70(2):343–365.
- Bengtsson, T., Bickel, P., and Li, B. (2008). Curse-of-dimensionality revisited: Collapse of the particle filter in very large scale systems. In Speed, T. and Nolan, D., editors, *Probability and Statistics: Essays in Honor of David A. Freedman*, pages 316–334. Institute of Mathematical Statistics, Beachwood, OH.
- Bickel, P. J., Klaassen, C. A. J., Ritov, Y., and Wellner, J. A. (1993). *Efficient and Adaptive Estimation for Semiparametric Models*. Johns Hopkins University Press, Baltimore.
- Bjørnstad, O. N. and Grenfell, B. T. (2001). Noisy clockwork: Time series analysis of population fluctuations in animals. *Science*, 293:638–643.

## References II

- Box, G. E. and Cox, D. R. (1964). An analysis of transformations. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)*, pages 211–252.
- Cox, D. and Snell, E. (1970). *The analysis of binary data*. Methuen and Co, London.
- Cox, D. R. (1972). Regression models and life-tables. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)*, 34:187–220.
- Doucet, A., Pitt, M. K., Deligiannidis, G., and Kohn, R. (2015). Efficient implementation of Markov chain Monte Carlo when using an unbiased likelihood estimator. *Biometrika*, 102:295–313.
- Efron, B. and Hinkley, D. V. (1978). Assessing the accuracy of the maximum likelihood estimator: Observed versus expected Fisher information. *Biometrika*, 65:457–483.

## References III

- Ellner, S. P., Bailey, B. A., Bobashev, G. V., Gallant, A. R., Grenfell, B. T., and Nychka, D. W. (1998). Noise and nonlinearity in measles epidemics: Combining mechanistic and statistical approaches to population modeling. *American Naturalist*, 151:425–440.
- Fasiolo, M., Pya, N., and Wood, S. N. (2016). A comparison of inferential methods for highly nonlinear state space models in ecology and epidemiology. *Statistical Science*, 31:96–118.
- Grenfell, B. T., Pybus, O. G., Gog, J. R., Wood, J. L. N., Daly, J. M., Mumford, J. A., and Holmes, E. C. (2004). Unifying the epidemiological and evolutionary dynamics of pathogens. *Science*, 303:327–332.
- Ionides, E. L. (2005). Maximum smoothed likelihood estimation. *Statistica Sinica*, 15:1003–1014.
- Ionides, E. L., Nguyen, D., Atchadé, Y., Stoev, S., and King, A. A. (2015). Inference for dynamic and latent variable models via iterated, perturbed Bayes maps. *Proceedings of the National Academy of Sciences of the USA*, 112:719–724.

## References IV

- King, A. A., Nguyen, D., and Ionides, E. L. (2016). Statistical inference for partially observed Markov processes via the R package pomp. *Journal of Statistical Software*, 69:1–43.
- Le Cam, L. (1986). *Asymptotic Methods in Statistical Decision Theory*. Springer, New York.
- Lindsay, B. G. and Li, B. (1997). On second-order optimality of the observed Fisher information. *Annals of Statistics*, 25(5):2172–2199.
- Murphy, S. A. and van der Vaart, A. W. (2000). On profile likelihood. *Journal of the American Statistical Association*, 95:449–465.
- Park, J. and Ionides, E. L. (2017). A guided intermediate resampling particle filter for inference on high dimensional systems. *Arxiv:1708.08543*.
- Rebeschini, P. and van Handel, R. (2015). Can local particle filters beat the curse of dimensionality? *The Annals of Applied Probability*, 25:2809–2866.

- Shrestha, S., King, A. A., and Rohani, P. (2011). Statistical inference for multi-pathogen systems. *PLoS Computational Biology*, 7:e1002135.
- Smith, R. A., Ionides, E. L., and King, A. A. (2017). Infectious disease dynamics inferred from genetic data via sequential Monte Carlo. *Molecular Biology and Evolution*, pre-published online, doi:10.1093/molbev/msx124.
- Wood, S. N. (2010). Statistical inference for noisy nonlinear ecological dynamic systems. *Nature*, 466:1102–1104.