


# What to do when exact Bayes is impossible?

## Some tools for approximate Bayesian inference

Umberto Picchini

picchini@chalmers.se

@uPicchini

<https://umbertopicchini.github.io>

Chalmers University of Technology and University of Gothenburg  
Sweden

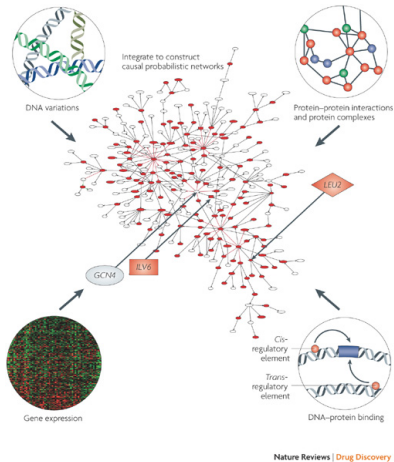
Notice, all text in [pink](#) (except this one) represents a URL that you can click.

- I will briefly introduce a methodology that has literally revolutionised statistical inference for complex models in the last 10-15 years.
- For the last 30 years advancements in computer hardware have enabled modellers to become more and more ambitious.
- Complex models are needed to make sense of advanced experiments and multivariate (large) datasets.

However the advancements of statistical algorithms didn't proceed at the same (fast) pace as hardware and modelling advancements.

**We wanted** to consider realistic model for our data, but often **we could not** because of the lack of flexible statistical methods.

Most real-life modelling is way more complex than examples from courses textbooks. The likelihood of the object below might be totally out of reach.



[Pic from Schadt et al. (2009) doi:10.1038/nrd2826]

Say that you have a complex network, with many edges, as in the previous figure.

You would like to run a Gibbs sampler...however you do not know how to sample from all the conditionals.

Ok, fear not, you can introduce a Metropolis-Hastings step to sample from those conditional that are unknown.

...oops, it turns out those conditionals have an unknown density. **We can't run Metropolis-within-Gibbs!**

More generally **we may not know the likelihood function explicitly** or it could be too computationally expensive to approximate.

And suppose the model is not of state-space type, where we would have ad-hoc trusted methods.

So what do we do?

What we typically want is the **likelihood function** for model parameters  $\theta$ :

- We have some data:  $\mathbf{y}^o$ .
- the likelihood function:  $p(\mathbf{y}^o|\theta)$
- We consider data as the outcome of some probabilistic model, and write  $\mathbf{y}^o \sim p(\mathbf{y}|\theta = \theta_0)$
- $\theta_0$  is the **unknown** *ground-truth* value of  $\theta$  that generated the data.

### Main issue

For realistically complex models, the likelihood function is **unavailable** in closed form.

Hence exact likelihood based inference is often **not possible**.

A paradigm shift is the concept of **generative model or simulator**.

You code a mathematical model  $\mathcal{M}(\boldsymbol{\theta})$  as an idealized representation of the phenomenon under study.

$$\boldsymbol{\theta}^* \rightarrow \mathcal{M}(\boldsymbol{\theta}^*) \rightarrow \mathbf{y}^*$$

As long as we are able to run an instance of the model, we simulate/generate artificial data  $\mathbf{y}^*$  with  $\mathbf{y}^* \sim p(\mathbf{y}^* | \boldsymbol{\theta} = \boldsymbol{\theta}^*)$ .

So we have obtained a random *realization*  $\mathbf{y}^*$  of the generative model  $\mathcal{M}(\boldsymbol{\theta})$

Therefore the simulator  $\mathcal{M}(\boldsymbol{\theta})$  defines the model pdf  $p(\mathbf{y} | \boldsymbol{\theta})$  **implicitly!**

## Trivial examples of generative models

**Example 1:** say that you get told that  $y \sim N(\mu, \sigma^2)$  but **assume you do not know how to sample from a general Gaussian** (just assume...).

However, given knowledge of how to sample  $u \sim N(0, 1)$ , then we also know how to **sample** from a generic Gaussian.

$$\mathcal{M} : y = \mu + \sigma \times u, \quad u \sim N(0, 1) \quad \text{independently}$$

then  $y \sim N(\mu, \sigma^2)$  with  $\theta = (\mu, \sigma)$ .

So the equation above is a **generative model**  $\mathcal{M}(\theta)$  for iid Gaussian draws. We only need some random input ( $u$ ) and then we can generate draws.



# Trivial examples of generative models

**Example 2:** stochastic Ricker model.

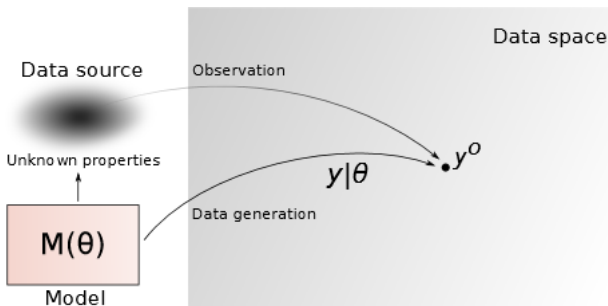
$$\mathcal{M}(\theta) = \begin{cases} \text{(observations): } y_t \sim_{\text{indep.}} \text{Poisson}(\phi N_t), & t = 1, \dots, T \\ \text{(unobservable process): } N_t = r \cdot N_{t-1} \cdot e^{-N_{t-1} + e_t}, & e_t \sim_{\text{iid}} \mathcal{N}(0, \sigma^2) \end{cases}$$

This **state-space model** can be used to describe the evolution in time of a population of size  $N_t$ . [We study this model later.](#)

Even though its likelihood is **analytically intractable** and given by

$$p(y_{1:T}|r, \phi, \sigma) = \int \prod_{t=1}^T (p(y_t|N_t, r, \sigma) p(N_t|N_{t-1})) \cdot dN_1 \cdots dN_T$$

we can still simulate the  $y_t \sim p(y|r, \phi, \sigma)$  from the (unknown) likelihood, because we can just run a computer model implementing  $\mathcal{M}(\theta)$  at any  $\theta = (r, \Phi, \sigma)$ .





We can use simulations from the generative model to produce inference about  $\theta$ , **without explicit knowledge of the likelihood**  $p(y|\theta)$ .

This is at the basis of **likelihood-free methods**

We are entering **simulation based inference**. Lots of literature available.

Some interesting discussions in **this review** (a bit biased towards the authors' own research):

## The frontier of simulation-based inference

Kyle Cranmer<sup>a,b,1</sup> , Johann Brehmer<sup>a,b</sup> , and Gilles Louppe<sup>c</sup>

<sup>a</sup>Center for Cosmology and Particle Physics, New York University, New York, NY 10003; <sup>b</sup>Center for Data Science, New York University, New York, NY 10011; and <sup>c</sup>Montefiore Institute, University of Liège, B-4000 Liège, Belgium

Edited by Jitendra Malik, University of California, Berkeley, CA, and approved April 10, 2020 (received for review November 4, 2019)

Many domains of science have developed complex simulations to describe phenomena of interest. While these simulations provide high-fidelity models, they are poorly suited for inference and lead to challenging inverse problems. We review the rapidly developing field of simulation-based inference and identify the forces giving additional momentum to the field. Finally, we describe how the frontier is expanding so that a broad audience can appreciate the profound influence these developments may have on science.

the simulator—is being recognized as a key idea to improve the sample efficiency of various inference methods. A third direction of research has stopped treating the simulator as a black box and focused on integrations that allow the inference engine to tap into the internal details of the simulator directly.

Amidst this ongoing revolution, the landscape of simulation-based inference is changing rapidly. In this review we aim to provide the reader with a high-level overview of the basic ideas behind both old and new inference techniques. Rather than

*One* of the possible simulation-based methods is ABC (approximate Bayesian Computation).

# ABC, approximate Bayesian computation

ABC is probably the most studied likelihood-free methodology.

The first and simplest ABC algorithm is **acceptance-rejection sampling**.

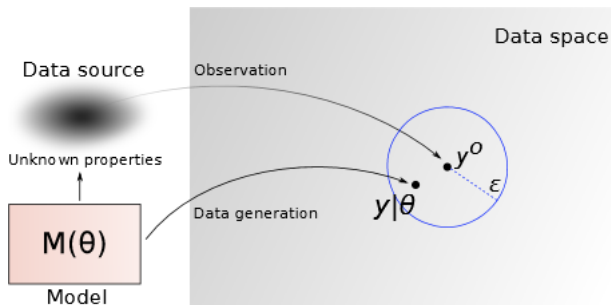
- 1 simulate from the prior  $\theta^* \sim \pi(\theta)$
- 2 plug  $\theta^* \rightarrow \mathcal{M}(\theta^*) \rightarrow \mathbf{y}^*$
- 3 if  $\|\mathbf{y}^* - \mathbf{y}^o\| < \epsilon$  accept  $\theta^*$  **otherwise discard**. Go to step 1 and repeat as many times as needed to obtain  $N$  accepted draws.

Each accepted pair  $(\theta^*, \mathbf{y}^*)$  is from the augmented-posterior  $\pi_\epsilon(\theta, \mathbf{y}^* | \mathbf{y}^o)$ .

But we do not really care for  $\mathbf{y}^*$ , so if we retain only accepted  $\theta^*$  then

$$\theta^* \sim \pi_\epsilon(\theta | \mathbf{y}^o)$$

**No likelihood was explicitly involved**, only implicitly via simulation!



Simulated data  $y^*$  inside the blue circle correspond to accepted parameters  $\theta^*$ .

## Which posterior are we targeting?

Acceptance-rejection sampling produces draws from the joint “augmented posterior”  $\pi_{\epsilon}(\boldsymbol{\theta}, \mathbf{y}^* | \mathbf{y}^o)$  where

$$\pi_{\epsilon}(\boldsymbol{\theta}, \mathbf{y}^* | \mathbf{y}^o) \propto \mathbb{I}_{\epsilon}(\mathbf{y}^*, \mathbf{y}^o) p(\mathbf{y}^* | \boldsymbol{\theta}^*) \pi(\boldsymbol{\theta}^*)$$

where  $\mathbb{I}_{\epsilon}(\mathbf{y}^*, \mathbf{y}^o)$  equals 1 if  $\|\mathbf{y}^* - \mathbf{y}^o\| \leq \epsilon$  and 0 otherwise.

However, in reality we do not need to store the  $\mathbf{y}^*$  (we can just discard those immediately after we have evaluated  $\|\mathbf{y}^* - \mathbf{y}^o\| \leq \epsilon$ ), and then  $\boldsymbol{\theta}^* \sim \pi_{\epsilon}(\boldsymbol{\theta} | \mathbf{y}^o)$  where

$$\pi_{\epsilon}(\boldsymbol{\theta} | \mathbf{y}^o) \propto \pi(\boldsymbol{\theta}^*) \int_{\mathcal{Y}} \mathbb{I}_{\epsilon}(\mathbf{y}^*, \mathbf{y}^o) p(\mathbf{y}^* | \boldsymbol{\theta}^*) d\mathbf{y}^*$$

Say that  $\theta^*$  has been accepted by the ABC acceptance-rejection sampling.

Then:

- if  $\epsilon = 0$  then  $\theta^* \sim \pi(\theta|y^o)$ , the **exact posterior**
- if  $\epsilon = \infty$  then  $\theta^* \sim \pi(\theta)$ , the **prior**

## Toy model

Let's try something really trivial. We show how acceptance-rejection can easily become inefficient.

The Weibull distribution has two positive parameters,  $a$  = “shape” and  $b$  = “scale”.

This is a **tractable distribution**. We use a tractable case study as it is pedagogically useful, since we can compare ABC inference with exact inference.

Its pdf for  $x > 0$  is

$$f(x) = (a/b)(x/b)^{a-1} \exp(-(x/b)^a)$$

and 0 otherwise.



Suppose we have  $n = 5$  i.i.d. observations  $y_i \sim \text{Weibull}(2, 5)$ .

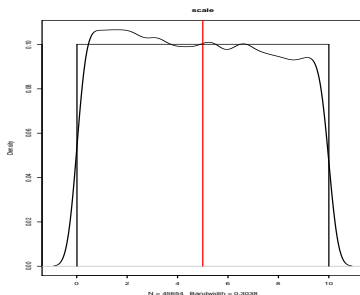
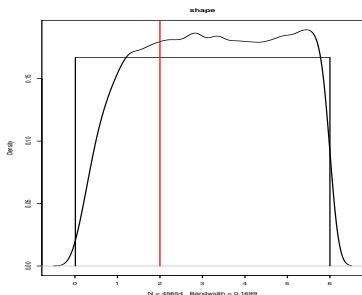
Want to estimate parameters of the Weibull, so  $\theta_0 = (2, 5) = (a, b)$  are the true values.

- take  $\| \mathbf{y}^o - \mathbf{y}^* \| = \sqrt{\sum_{i=1}^n (y_i^o - y_i^*)^2}$  (you can try a different distance, this is not really crucial).
- We'll use different thresholds  $\epsilon$ .
- Run 50,000 iterations of acceptance-rejection.

Notice, in this case we prefix the number of iterations. So we do not know in advance how many accepted parameters we get. Instead in the algorithm we previously defined, I wrote that we keep repeating until a number  $N$  of acceptances is obtained.

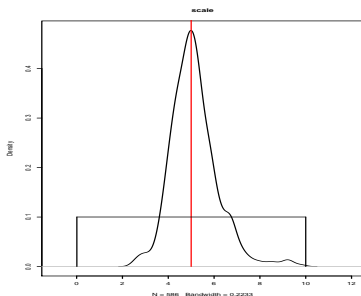
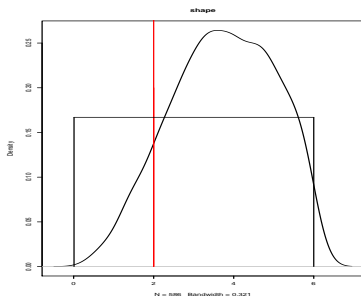
Wide priors for the “shape” parameter  $a \sim U(0.01, 6)$  and “scale”  $b \sim U(0.01, 10)$ .

Try  $\epsilon = 20$ . True parameter values in red.



We are evidently sampling from the prior. Must reduce  $\epsilon$ . About 92% draws were accepted. Way too large percentage!

Reduce  $\epsilon$  from  $\epsilon = 20$  to  $\epsilon = 3$



About 1% of the produced simulations has been accepted.

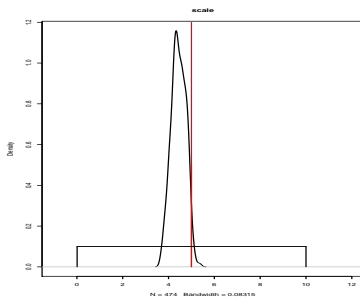
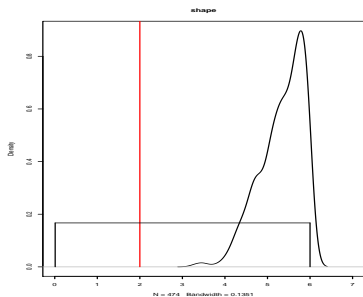
Of course  $n = 5$  is a very small sample size, so inference quality is necessarily limited, but you got an idea of the ABC method.

An acceptance rate of about 1% is often implemented in many studies as a good tradeoff between computational effort and statistical accuracy.

# Curse of dimensionality

- results will degrade for a larger sample size  $n$  because of a “necessarily too large”  $\epsilon$ ;
- even for a moderately long dataset  $\mathbf{y}^o$ , how likely is that we simulate a  $\mathbf{y}^*$  such that  $\sum_{i=1}^n (y_i - y_i^*)^2 < \epsilon$  for **small**  $\epsilon$ ?  
Very unlikely.
- inevitably, we'll be forced to enlarge  $\epsilon$  thus degrading the quality of the inference.
- Serious trade-off between computational efficiency and statistical precision.

Here we take  $n = 200$ . To compare with our “best” previous result, we use  $\epsilon = 31$  (to obtain again a 1% acceptance rate on 50,000 iterations).



Notice shape is completely off!

The approach is just **not going to be of any practical use with large datasets.**

# Break the curse of dimensionality

**Compress** data information using some **summary statistics**  $S(\mathbf{y})$ .

Example:  $S(\mathbf{y})$  may contain sample mean, standard deviation, autocorrelations, quantiles etc.

**Idea:** instead of comparing  $\mathbf{y}^o$  with  $\mathbf{y}^*$ , compare  $S(\mathbf{y}^o)$  with  $S(\mathbf{y}^*)$ .

Requirements:

- $S(\cdot)$  should be “informative” regarding  $\theta$ , as we give up on using the full data  $\mathbf{y}$ .
- $S(\cdot)$  should not be too large. Ideally  $\dim(\mathbf{S}) \equiv \dim(\theta)$  [Fearnhead & Prangle '12].<sup>1</sup>

---

<sup>1</sup>Fearnhead & Prangle (2012). JRSS-B, 74(3), 419-474.

## Acceptance-rejection with summaries (Pritchard et al.<sup>2</sup>)

- 1 simulate from the prior  $\theta^* \sim \pi(\theta)$
- 2 simulate  $\mathcal{M}(\theta^*) \rightarrow \mathbf{y}^*$ , compute  $S(\mathbf{y}^*)$
- 3 if  $\| S(\mathbf{y}^*) - S(\mathbf{y}^o) \| < \epsilon$  store  $\theta^*$ . Go to step 1 and repeat as many times as needed to obtain  $N$  accepted draws.

Samples are from  $\pi_\epsilon(\theta|S(\mathbf{y}^o))$

with

$$\pi_\epsilon(\theta|S(\mathbf{y}^o)) \propto \pi(\theta^*) \int_{\mathcal{Y}} \mathbb{I}_{A_{\epsilon, \mathbf{y}^o}}(\mathbf{y}^*) p(\mathbf{y}^*|\theta^*) d\mathbf{y}^*$$

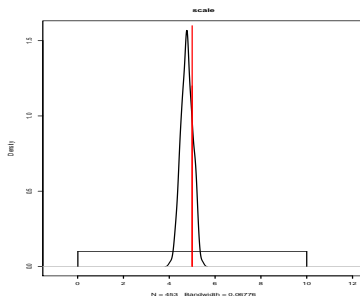
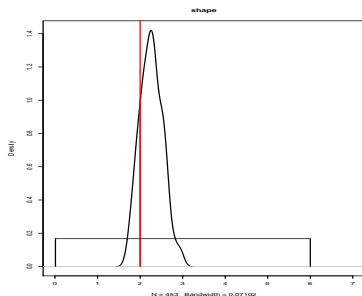
$$A_{\epsilon, \mathbf{y}^o}(\mathbf{y}^*) = \{\mathbf{y}^* \in \mathcal{Y}; \| S(\mathbf{y}^*) - S(\mathbf{y}^o) \| < \epsilon\}.$$

---

<sup>2</sup>Pritchard et al. 1999, Molecular Biology and Evolution, 16:1791-1798.

## Weibull example, reprise with $n = 200$

Set  $S(\mathbf{y}) = (\text{sample mean of } \mathbf{y}, \text{sample SD of } \mathbf{y})$ . Set  $n = 200$  observations.  
Use  $\epsilon = 0.35$ .



This time we have captured **both shape and scale** (with 1% acceptance).

Also, enlarging  $n$  would not cause problems, thanks to  $S(\cdot)$ .



Using summary statistics clearly introduces a further level of approximation. Except when  $S(\cdot)$  is *sufficient* for  $\theta$  (carries the same info about  $\theta$  as the whole  $\mathbf{y}^o$ ).

When  $S(\cdot)$  is a set of sufficient statistics for  $\theta$ ,

$$\pi_{\epsilon}(\theta|S(\mathbf{y}^o)) = \pi_{\epsilon}(\theta|\mathbf{y}^o)$$

However when the distribution of  $\mathbf{y}^o$  is not in the *exponential family*, we basically have **no hope to construct sufficient statistics**.

For the sake of completeness: we could achieve exact inference under the rather stringent hypotheses that

- 1  $S$  is sufficient, and simultaneously that
- 2  $\epsilon = 0$

Note, (2) is **completely impractical** unless  $\mathbf{y} \in \mathcal{Y}$ , with  $\mathcal{Y}$  a discrete set having few possible values.

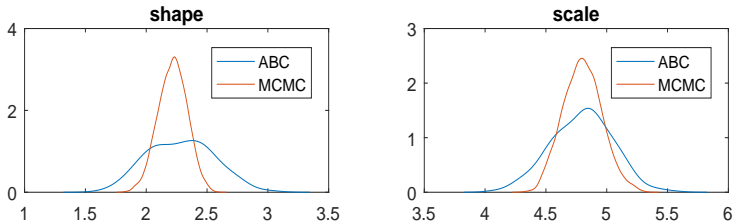
ABC methodology is only apparently simple.

A large amount of research has been produced in the last 20 years to improve over the basic algorithm.

We will see an example soon, but first...

For this toy model, exact inference is possible.

**Did we obtain an accurate approximation to the exact posterior?**



Exact posteriors (via MCMC) are in red.

ABC has 2 (+1) sources of approximation:

- we used arbitrary **non-sufficient statistics**.
- $\epsilon > 0$
- Monte Carlo approximation due to finite samples (but that's obvious)

The apparent simplicity of ABC acceptance-rejection should not promote “lazy science”, just because it can be easily run.

At any time you should strive to use (or search for) the methods that allow you to implement exact inference.

ABC should be a last-resort if better methods are unavailable, since it is unknown how much approximate the ABC results will be (you typically cannot compare with exact inference, unlike with the previous toy model).

## Beyond ABC rejection

ABC rejection is the simplest example of ABC algorithm.

It generates independent draws and can be coded into an **embarrassingly parallel algorithm**. However it can be very inefficient.

Parameters are proposed from the prior  $\pi(\theta)$ . **A prior does not exploit the information of already accepted parameters.**

Unless  $\pi(\theta)$  is similar to the posterior, many proposals will be rejected for moderately small  $\epsilon$ . This **worsen with increasing dimension of  $\theta$** .

A natural approach is to consider ABC within an MCMC algorithm.

In a MCMC with random walk proposals, the proposed parameter explores a neighbourhood of the last accepted parameter.

# ABC-MCMC [Marjoram et al. 2003]

Integrating ABC within MCMC is **very simple** [Marjoram et al. 2003]

Notation: write  $\mathbf{s}^* \equiv S(\mathbf{y}^*)$ ,  $\mathbf{s}^o \equiv S(\mathbf{y}^o)$ .

and  $\mathbb{I}_\epsilon(\mathbf{s}^*, \mathbf{s}^o)$  equals 1 if  $\|\mathbf{s}^* - \mathbf{s}^o\| < \epsilon$ , and 0 otherwise.

- 1 sample proposal  $\theta^* \sim q(\theta^* | \theta^\#)$
- 2 plug  $\theta^* \rightarrow \mathcal{M}(\theta^*) \rightarrow \mathbf{y}^* \rightarrow \mathbf{s}^*$
- 3 compute acceptance ratio:

$$\text{ratio} := \frac{\mathbb{I}_\epsilon(\mathbf{s}^*, \mathbf{s}^o) \pi(\theta^*)}{\mathbb{I}_\epsilon(\mathbf{s}^\#, \mathbf{s}^o) \pi(\theta^\#)} \times \frac{q(\theta^\# | \theta^*)}{q(\theta^* | \theta^\#)}$$

- 4 simulate  $u \sim U(0, 1)$ , accept  $\theta^*$  if  $u < \text{ratio}$ .

The previous algorithm produces dependent samples from the “augmented” posterior  $\pi_{\epsilon}(\theta, \mathbf{s}^* | \mathbf{s}^o)$ .

This means that if we disregard  $\mathbf{s}^*$ , and retain only  $\theta^*$ , we have

$$\theta^* \sim \pi_{\epsilon}(\theta | \mathbf{s}^o)$$

This is just another way to sample from an approximated posterior. Using a more informed proposal function than the prior.

## Example: stochastic Ricker model

$$\begin{cases} \text{(observations): } y_t \sim \text{Poisson}(\phi N_t), & t = 1, \dots, T \\ \text{(unobservable process): } N_t = r \cdot N_{t-1} \cdot e^{-N_{t-1} + e_t}, & e_t \sim_{iid} \mathcal{N}(0, \sigma^2) \end{cases}$$

It can be used to describe the evolution in time of a population of size  $N_t$ .

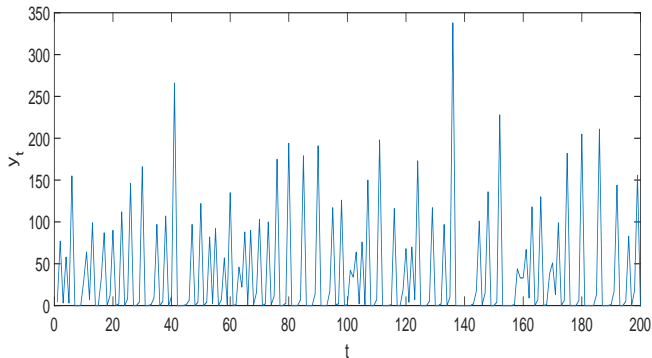
- $r$  is the intrinsic growth rate of the population;
- $\phi$  is a scale parameter
- $e_t$  interpreted as environmental noise.

This is a state-space model, as the dynamics of  $\{N_t\}$  are Markovian and we assume measurements  $y_{1:T}$  to be conditionally independent given  $\{N_t\}$ .



## the data

We simulated 200 time points from the model, with  
 $\log r = 3.8$ ,  $\log \phi = 2.3$ ,  $\log \sigma = -1.2$



## Summary statistics

We used the 13 summary statistics suggested in [Wood 2010](#). These include:

- the sample mean of observations  $\bar{y}$ ;
- number of zeros in the dataset;
- autocovariances up to lag 5;
- and six more summaries...(not important to be mentioned here, see the reference above or the provided MATLAB code).

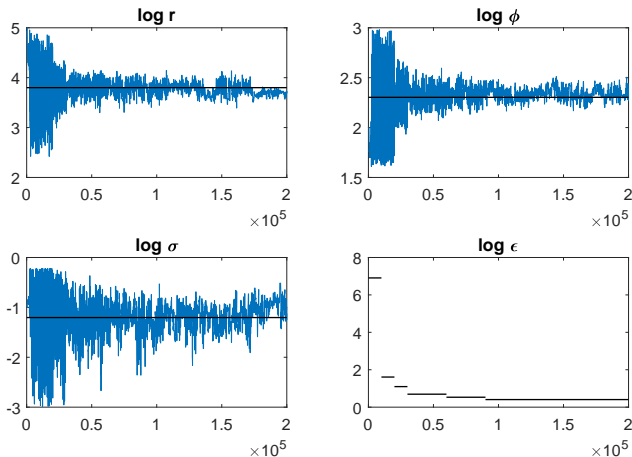
So we have  $s(\mathbf{y}^o) = (\bar{y}, \#zeros, autocov \text{ lag}1, \dots, autocov \text{ lag}5, \dots)$ .

Priors:

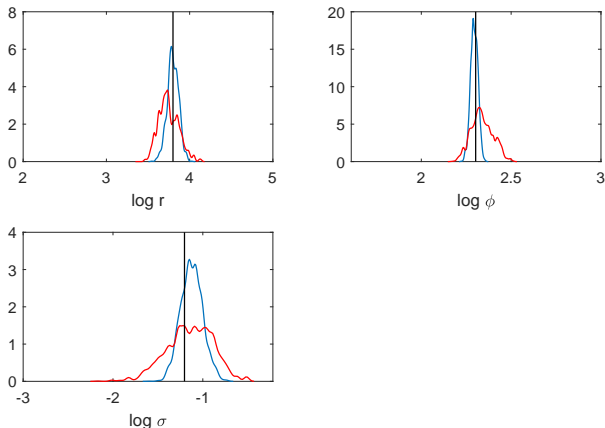
- $\log r \sim U(2, 5)$
- $\log \Phi \sim U(1.61, 3)$
- $\log \sigma \sim U(-3, -0.22)$

# ABC-MCMC traces

We performed 200,000 ABC-MCMC iterations with decreasing  $\epsilon_t$ .



ABC-MCMC is in red.



In blue is inference via “particle MCMC”. No time to talk about this but the blue one is basically exact inference (up to Monte Carlo error).

## Preliminary conclusions

- ABC allows you to produce approximate inference for models having an intractable/unknown likelihood function.
- in our examples **we never needed to know the likelihood**;
- ABC only requires the ability to simulate artificial data, but is not a silver bullet.
- the main difficulty is how to specify summary statistics that are “informative” for  $\theta$ .
- when summaries are not informative and  $\epsilon$  is too large results will be poor.
- tuning ABC is not straightforward. Many available resources though, see the final slides...

## Other likelihood-free methods

Likelihood-free methods date back to at least [Diggle and Gratton \(1984\)](#) and [Rubin \(1984, p. 1160\)](#)

More recent examples:

- **Indirect Inference** ([Gourieroux and Ronchetti 1993](#));
- for state-space models (Markov processes observed with noise) the **bootstrap filter** of [Gordon, Salmond and Smith \(1993\)](#)
- **Synthetic Likelihoods** method of [Wood \(2010\)](#)
- Lots of more recent machine learning literature: see [this review](#).

# Are approximations any worth?

Why should we care about approximate methods?

Well, we know the most obvious answer: it's because this is what we do when exact methods are impractical. No big news...

But I am more interested in the following phenomenon, which I noticed by direct experience:

- Many scientists seem to get intellectual fulfilment by using exact methods, leading to exact inference.
- What we might not see is when they fail to communicate that they (consciously or unconsciously) pushed themselves to formulate **simpler models (too simple?!), so that exact inference could be achieved.**

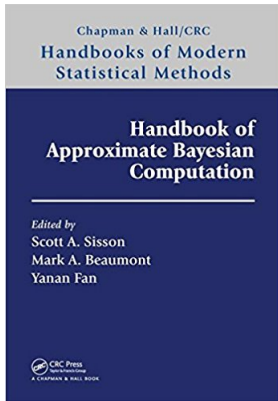
## Gelman and Rubin, 1996

“[...] as emphasized in Rubin (1984), one of the great scientific advantages of simulation analysis of Bayesian methods is the freedom it gives the researcher to formulate appropriate models **rather than be overly interested in analytically neat but scientifically inappropriate models.**”



The ABC methodology is quite mature but still evolving.  
Comprehensive 2018 monography (many chapters available on arxiv):

Sisson, Fan, Beaumont. (2018). Handbook of approximate Bayesian computation. Chapman and Hall/CRC.



## Blog posts and slides (coloured links are clickable)

- 1 [Christian P. Robert](#) often blogs about ABC (and beyond: it's a fantastic blog!)
- 2 an [intro to ABC](#) by Darren J. Wilkinson
- 3 Two blog posts by Rasmus Bååth [here](#) and [here](#)
- 4 Tons of slides at [Slideshare](#).

## Software (coloured links are clickable)

- [EasyABC](#), R package. Research [article](#).
- [abc](#), R package. Research [article](#)
- [abctools](#), R package. Research [article](#). Focusses on tuning.
- [pyABC](#), Python package.
- [ABCpy](#), Python package.
- A list with more options [here](#) .
- [examples](#) with implemented model simulators (useful to incorporate in your programs).

# Reviews

Accessible reviews:

- 1 [Sunnåker et al. 2013](#)
- 2 (with applications in ecology) [Hartig et al. 2013](#)

Fairly extensive reviews:

- 1 [Sisson and Fan 2010](#)
- 2 (with applications in ecology) [Beaumont 2010](#)
- 3 [Marin et al. 2010](#)

Review specific for dynamical models:

- 1 [Jasra 2015](#)

# Determination of summary statistics

- 1 review paper by [Blum et al. 2013](#) on dimension reduction methods for ABC;
- 2 [Fearnhed and Prangle 2012](#) (a JRSS-B discussion paper).
- 3 [Wiqvist et al. 2019](#) using deep learning.