# Simulation-based Inference %>%

**Concept and ELFI - tutorial %>%**
**ProbAI2022**

**Henri Pesonen**
**Oslo Center for Biostatistics and Epidemiology**

@henri_pesonen

# Concepts

- Bulding blocks of Likelihood-free inference

- Sampling based LFI methods

- Surrogate based methods

- Active learning

**Tutorial available as a notebook in Google Colab**

- Basics of model building and exploration

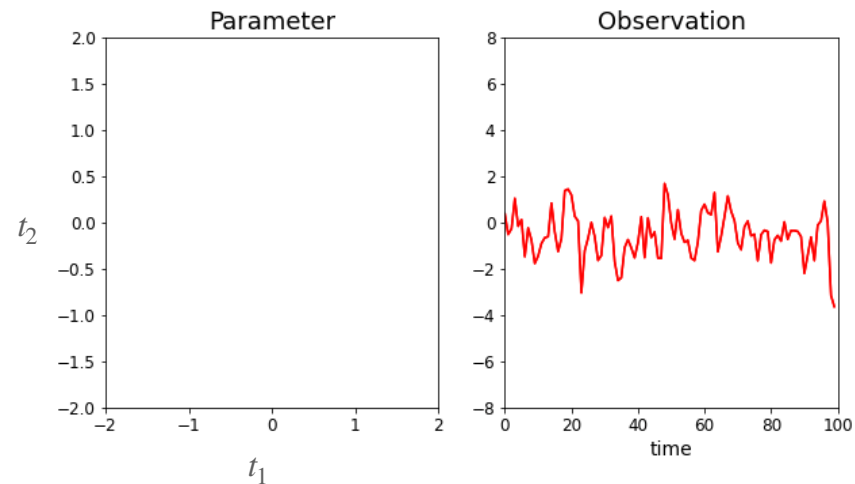- Choosing and using inference methods

# Simulation-based Inference



- The phenomena of the world are frequently investigated *in silico* - i.e. using complex computer simulations with great detail

- The simulators are controlled by a set of parameters that we want to infer based on the observations collected of the phenomena

- Complexity of the simulators often prohibits the access to the most important tool of statistical inference - likelihood function

# Example

## MA(2) model

- Simple time series model

- $x_t = w_t + t_1 w_{t-1} + t_2 w_{i-2}, \quad w_t \sim \text{Normal}(0,1)$

# Example
## Transmissions of bacterial infections in daycare centers.

- Cross-sectional data from a stochastic SIS-model

- Continuous-time Markov process with transition probabilities:

$$P(I_{is}(t + dt) = 1 \mid I_{is}(t) = 0) = \theta_1 \cdot E_s(I(t)) + \theta_2 \cdot P_s, \quad \text{if} \quad I_{i1}(t) + \cdots + I_{iN_s}(t) = 0$$

$$P(I_{is}(t + dt) = 1 \mid I_{is}(t) = 0) = \theta_3 \cdot (\theta_1 \cdot E_s(I(t)) + \theta_2 \cdot P_s), \quad \text{otherwise}$$

$$P(I_{is}(t + dt) = 0 \mid I_{is}(t) = 1) = \gamma$$

- $I_{is}(t)$ is the status of carriage of strain s for individual i.

- $E_s(I(t))$ is the probability of sampling the strain s

- $\theta_1$ is the rate of transmission from other children at the DCC

- $\theta_2$ is the rate of transmission from the community outside the DCC

- $\theta_3$ scales the rate of an infected child being infected with another strain

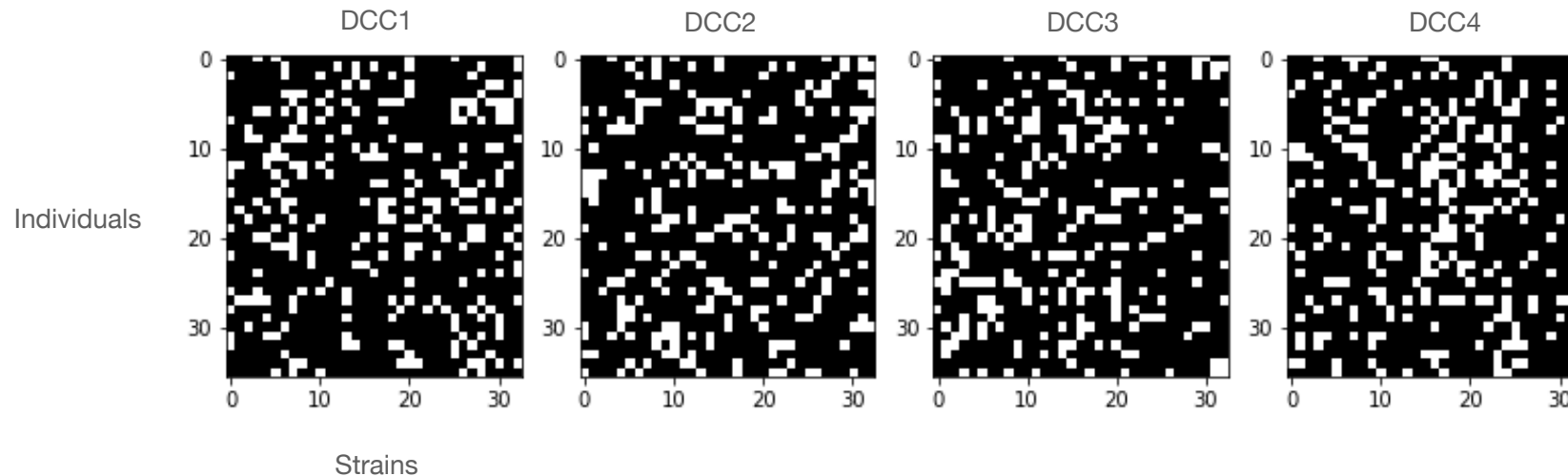- $\gamma$ is the relative probability of healing from a strain (scaled to 1)

# Example

## Transmissions of bacterial infections in daycare centers.

$$P(I_{is}(t + dt) = 1 \,|\, I_{is}(t) = 0) = \theta_1 \cdot E_s(I(t)) + \theta_2 \cdot P_s, \quad \text{if} \quad I_{i1}(t) + \cdots + I_{iN_s}(t) = 0$$

$$P(I_{is}(t + dt) = 1 \,|\, I_{is}(t) = 0) = \theta_3 \cdot (\theta_1 \cdot E_s(I(t)) + \theta_2 \cdot P_s), \quad \text{otherwise}$$
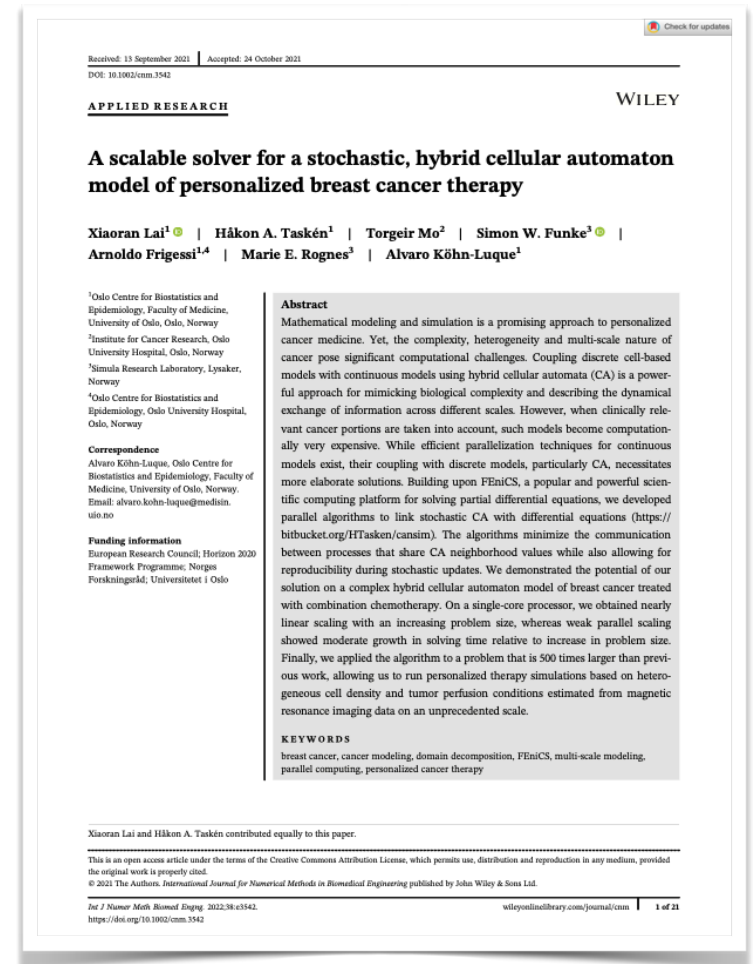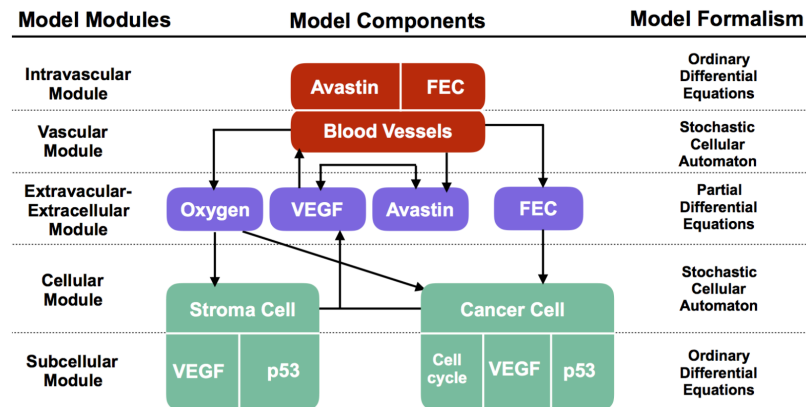
$$P(I_{is}(t + dt) = 0 \,|\, I_{is}(t) = 1) = \gamma$$

# Example
## Personalised medicine

- Model for evolution of breast cancer treated with combination chemotherapy

- Describes the evolution of cancer cells, blood vessels, Oxygen, VEGF and Avastin

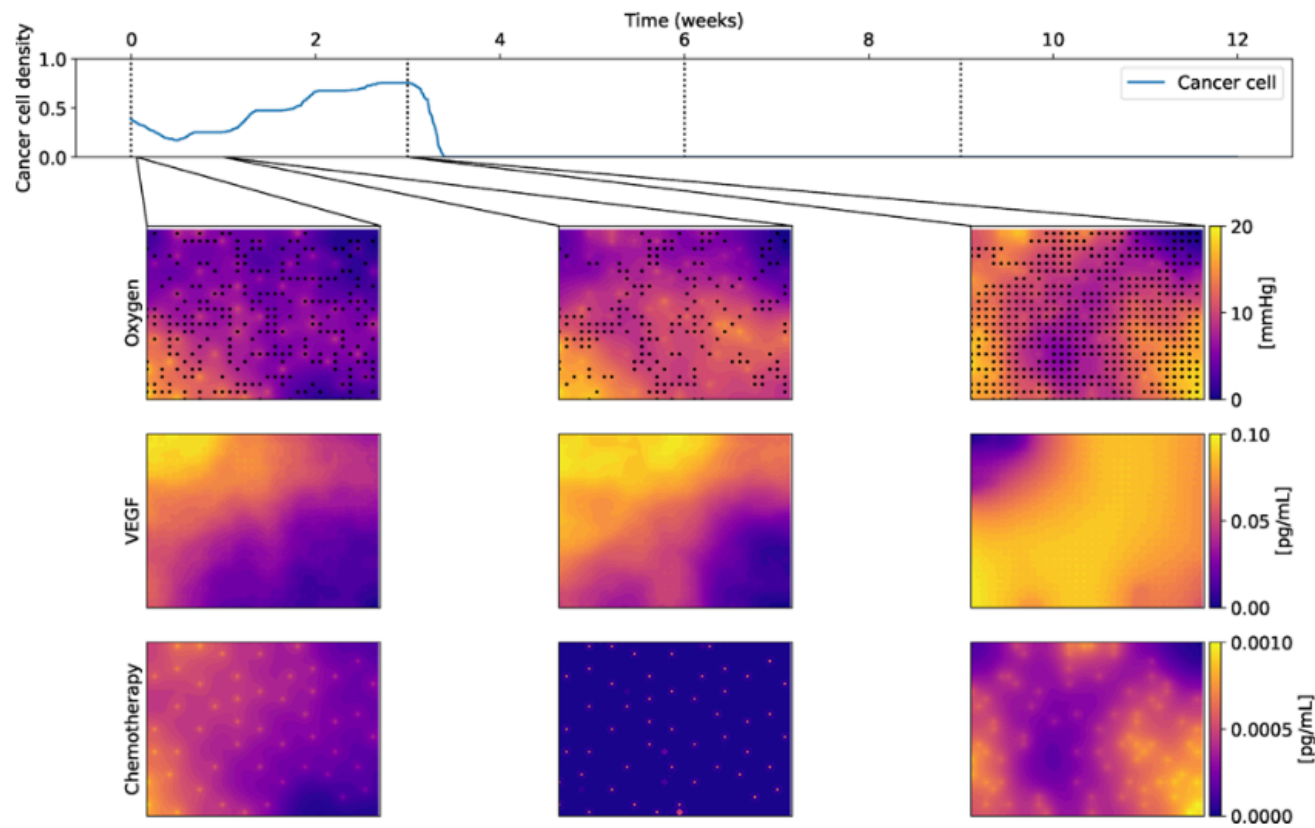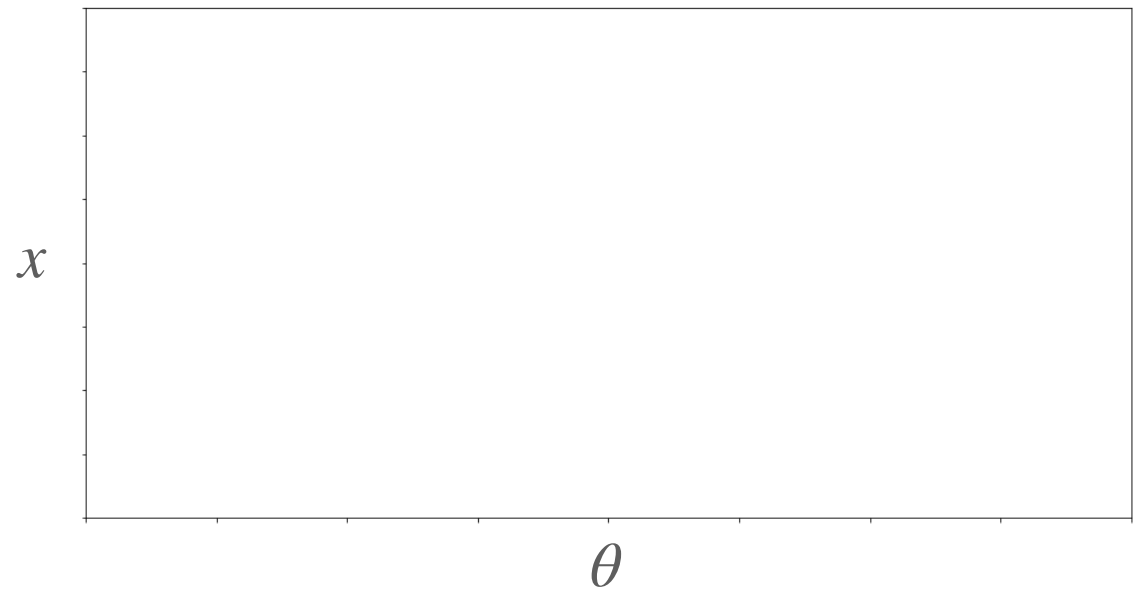# Example
## Personalised medicine



Figure from Lai et al, 2021

# Simulators %>% Likelihood-free Inference

# Simulator

- A computer program defined as $x \sim p(x \mid \theta)$ that has

  - input parameters $\theta$

  - stochastic output $x$
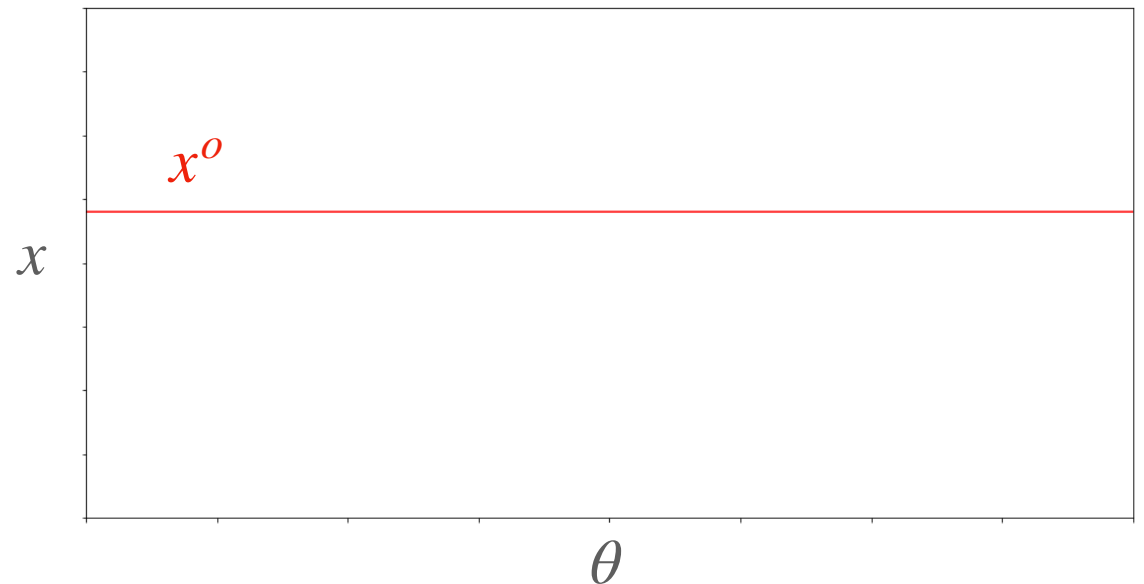
# Simulator

- The data produced by the simulator can be basically in any format

- It can be e.g.

  - Single time series

  - Set of independent data

  - Images

  - Distribution of data points

# Inference

- Observe data and infer the values of the parameters that generated them

  - Often based on likelihood $p(x^o \mid \theta)$

  - Bayesian approach $p(\theta \mid x^o) \propto p(x^o \mid \theta) p(\theta)$

  - Maximum likelihood $\arg\max_\theta p(x^o \mid \theta)$

- When data generating process (simulator) is defined as a set of rules to draw $x \sim p(x \mid \theta)$ it is often infeasible to formulate the analytical likelihood $p(x^o \mid \theta)$

# Inference without likelihood

- Use the capability to draw simulated data conditioned on the input parameters

- Likely true parameter values are thought to produce data that are *similar* to the observed data

# Inference without likelihood

- Use the capability to draw simulated data conditioned on the input parameters

- Likely true parameter values are thought to produce data that are *similar enough* to the observed data

- Approximate the posterior distribution as

$$p(\theta \mid x^o) \propto p(x^o \mid \theta)p(\theta) = \int \mathbb{1}_{\Omega(x^o)}(x)p(x \mid \theta)\mathrm{d}x\, p(\theta)$$

- Region $\Omega(x^o) = \{x : d(x, x^o) \leq \epsilon\}$ contains data that are similar enough to the realised observation

# Distance metric

- Acceptation region is defined by the distance metric

- Reasonable distance metrics depend on the data format, e.g.

  - Euclidean distance for low dimensional numerical outputs

  - L1 distance for data containing outliers

  - Quantiles or Wasserstein distance for distributions

# Data dimensionality

- Dimensionality of the data often causes problem when running LFI methods

  - In high dimensions it becomes increasing improbable to generate data close to the observed data

- The (current) standard approach is still to use summary statistics $S(\cdot)$

  - E.g $d(x, x^o) \approx d(S(x), S(x^o)) \quad (\approx \rho(x, x^o))$

- Sufficient statistics usually do not exist

- How to choose them?

# Selecting summary statistics

- An open problem

- Often we use bespoke summary statistics

  - Use **domain expertise** if available

  - **Explore** the simulator prior to inference

  - **Diagnose** the inference results

- Automatic algorithms for **selecting/constructing** the summary statistics

Summary statistics %>%
Distance function %>%
Threshold

# Rejection Approximate Bayesian Computation
## (ABC)

For $i = 1, \ldots, N$

$\quad \theta* \sim p(\theta)$

$\quad x* \sim p(x \mid \theta*)$

$\quad$ while $d(x*, x^o) > \epsilon$ :

$\quad\quad \theta* \sim p(\theta)$

$\quad\quad x* \sim p(x \mid \theta*)$

$\quad$ set $\theta_i = \theta*$

# Rejection ABC
**Alternative approach**

- Instead of choosing a fixed threshold we can instead sample a very large artificial data set and choose a fraction of samples that are most similar to the observed data

- Threshold can then be calculated after the simulation as the largest distance that was still accepted

- The approach can be used in smaller scale to find out reasonable threshold levels for sampling based ABC methods

# Rejection ABC uses samples from the prior

- Unless we have plenty of prior information about the parameters, sampling from prior is hardly effective

- Sequentially formulate importance sampling distributions with more probability mass in interesting regions

  - Sequential Monte Carlo ABC

# SMC-ABC

- Round 1 : SMC-ABC is initialised as rejection ABC with a loose threshold $\epsilon_1$

  - $\theta_i^{(1)} \sim p(\theta \mid d(S(x), S(x^o)) \leq \epsilon_1), \quad i = 1, \ldots, N$

  - Set weight $w_i^{(1)} = N^{-1}$

  - Calculate sample variances (for each dimension $j = 1, \ldots, M$)

  $$\hat{\mu} = \sum_{i=1}^{N} \frac{\theta_i}{N}, \quad \hat{\sigma}_j^2 = \sum_{i=1}^{N} \frac{1}{N}(\theta_{i_j} - \hat{\mu}_{i_j})^2$$

- Set proposal density $q(\theta^{(2)} \mid \theta^{(1)}) = \text{Normal}(\theta^{(1)}, 2 \cdot \text{diag}(\hat{\sigma}_1^2, \ldots, \hat{\sigma}_M^2))$

# SMC-ABC

- Rounds $t = 2, \ldots, T$: Set a tighter threshold $\epsilon_t < \epsilon_{t-1}$

  - (1) Select $\theta_i^{(t-1)}$ with probability $\propto w_i^{(t-1)}$

  - (2) Draw $\theta* \sim \text{Normal}(\theta_i^{(t-1)}, \text{diag}(\hat{\sigma}_1^2, \ldots, \hat{\sigma}_M^2))$

  - (3) Simulate $x* \sim p(x \mid \theta*)$

  - Repeat (1)-(3) until $d(S(x*), S(x^o)) < \epsilon_t$ and set $\theta_i^{(t)} = \theta*$

  - Set weights $w_i^{(t)} \propto p(\theta_i^{(t)}) \cdot \left[ \sum_{k=1}^{N} w_k^{(t-1)} \sum_{j=1}^{M} \phi \left( \frac{\theta_{i_j}^{(t)} - \theta_{k_j}^{(t)}}{\hat{\sigma}_j^{(t-1)}} \right) \right]^{-1}$

  - Calculate weighted sample variances $\hat{\sigma}_j^2$ (for each dimension $j = 1, \ldots, M$)

  - Set proposal density $q(\theta^{(t+1)} \mid \theta^{(t)}) = \text{Normal}(\theta^{(t)}, 2 \cdot \text{diag}(\hat{\sigma}_1^2, \ldots, \hat{\sigma}_M^2))$

Summary statistics %>%
Distance function %>%
Threshold %>%
Things that make you go *hmmm?*

# Issues to be aware of

- Summary statistics may not catch relevant features of the data

    - can decrease dimension too much - lose information

    - don't decrease dimension enough - didn't solve the problem

    - can be correlated - redundant dimensions

    - etc

- How multiple summary statistics/elements of a summary statistic contribute to the distance?

    - E.g. Multiple summary statistics can each have wildly different scales

# Reference

Chapman & Hall/CRC
Handbooks of Modern
Statistical Methods

**Handbook of
Approximate Bayesian
Computation**

Edited by
Scott A. Sisson
Yanan Fan
Mark A. Beaumont

CRC Press
A CHAPMAN & HALL BOOK

- For those who want to read more

# Surrogate models

- Alternative approach to likelihood-free inference is to construct surrogate models for various parts of the system

  - Synthetic likelihood is one of the first surrogate methods

  - At $\theta$ approximate
  $$p(x \mid \theta) \approx \text{Normal}(x \mid \hat{\mu}_N(\theta), \hat{\Sigma}_N(\theta)) \text{ with an sample}$$
  of size $N$ drawn from the model



Figure 1 | Measuring fit of the Ricker model. a, Population data simulated from the Ricker model in the text, observed under Poisson sampling ($\log(r) = 3.8$, $\sigma = 0.3$, $\phi = 10$). b, The log joint probability density, $\log(f(y, \mathbf{e}))$, of data, y, and random process noise terms, $\mathbf{e}$, plotted against the value of the first process noise deviate, $e_1$, with the rest of $\mathbf{e}$ and y held fixed. c, $\text{Log}(f(y, \mathbf{e}))$ plotted against model parameter r, again with $\mathbf{e}$ and y held fixed. d, The log synthetic likelihood, $l_s$, plotted against $\log(r)$ for the Ricker model and the data given in a ($N_s = 500$).

# Example

- Construct approximate likelihood at an arbitrary parameter value

Gaussian
empirical estimate

Set of samples
at $\theta$

$\approx p(\theta \mid x^o)$



$x^o$

$x$

$\theta$

$\theta$

$\theta$

# Synthetic likelihood is hardly efficient
**Surrogate is fitted at each parameter value separately**

- BOLFI - Bayesian optimization for likelihood-free inference:

  - Model the discrepancy as a function of the parameter $\Delta(\theta) = d(x(\theta), x^o)$ conditioned on simulated data $\{(\theta_i, \Delta(\theta_i))\}_{i=1}^t$ using a Gaussian process

$$\Delta(\theta) \mid \{(\theta_i, \Delta(\theta_i))\}_{i=1}^t \sim \text{GP}(\mu_{1:t}(\theta), v_{1:t}(\theta) + \sigma_n^2)$$

# Synthetic likelihood is hardly efficient
## Surrogate is fitted at each parameter value separately

- BOLFI - Bayesian optimization for likelihood-free inference:

  - Likelihood can be approximated from the surrogate (pointwise) by

$$p(x^o \mid \theta) \approx \Phi\left(\frac{\epsilon - \mu_{1:t}(\theta)}{\sqrt{v_{1:t}(\theta) + \sigma_n^2}}\right)$$



Figure from M. Järvenpää, "Efficient Acquisition Rules for Model-Based Approximate Bayesian Computation", 2019

# GP surrogates can utilise active learning

- Different strategies for selecting parameter values where to query the simulator

  - Reduce the number queries required to produce reasonable approximations to posterior/likelihood

- Usually based on optimization

  - Parameters that produce minimum discrepancies

  - Parameters that decrease most the uncertainty about the posterior

# BOLFI

# BOLFI

- Find parameter values that minimize discrepancy function

  - Black-box optimization

  - Acquisition strategies balance exploration and exploitation

# BOLFI

- Find parameter values that minimize discrepancy function

  - Black-box optimization

  - Acquisition strategies balance exploration and exploitation

$$p(x^o \mid \theta) \approx \Phi \left( \frac{\epsilon - \mu_{1:t}(\theta)}{\sqrt{v_{1:t}(\theta) + \sigma_n^2}} \right)$$

# Minimizing distance is often not optimal
## Ultimate goal is to approximate posterior distribution

- How to choose query points that are most informative about the posterior distribution

- Active learning strategies can be based on more reasonable goals

# Acquisition functions

- Lower Confidence Bound Selection Criterion

- Maximum Variance

- Randomized Maximum Variance

- Expected Integrated Variance

# LCBSC
## Lower Confidence Bound Selection Criteria for minimizing the distance

- Selecting the query point for round $t$ is a two part process. First optimise LCBSC

$$\theta* = \arg\min_{\theta} \mu_{1:t}(\theta) - \sqrt{\eta_t^2 v_{1:t}(\theta)} \quad , \quad \eta_t^2 = 2 \cdot \log\left(\frac{t^{2 \cdot d + 2}\pi^2}{3 \cdot \xi_\eta}\right)$$

- Then sample the next query point from truncated Gaussian

$$\theta_{t+1} \sim \mathsf{TN}(\theta*, \Sigma_{\mathsf{acq}}, \Omega)$$

- where $\Sigma_{\mathsf{acq}}$ and $\Omega$ are a tunable parameter balancing exploration/explotation and the optimization region, respectively

# MaxVar
**The maximum variance acquisition method**

- The next evaluation point is acquired where the variance of the unnormalised approximate posterior is maximised

$$\theta_{t+1} = \arg\max_{\theta} \text{Var}(p(\theta) \cdot p_a(\theta))$$

$$p_a(\theta) = \Phi\left(\frac{\epsilon - \mu_{1:t}(\theta)}{\sqrt{v_{1:t}(\theta) + \sigma_n^2}}\right)$$

- $\epsilon$ is the ABC threshold, $\mu_{1:t}$ and $v_{1:t}$ are determined by the GP surrogate, $\sigma_n^2$ is the noise.

# RandMaxVar
**The randomized maximum variance acquisition method**

- The next evaluation point is drawn randomly from the density corresponding to the variance of the posterior

$$\theta_{t+1} \sim q(\theta), \text{ where } q(\theta) \propto \text{Var}(p(\theta) \cdot p_a(\theta))$$

$$p_a(\theta) = \Phi \left( \frac{\epsilon - \mu_{1:t}(\theta)}{\sqrt{v_{1:t}(\theta) + \sigma_n^2}} \right)$$

- $\epsilon$ is the ABC threshold, $\mu_{1:t}$ and $v_{1:t}$ are determined by the GP surrogate, $\sigma_n^2$ is the noise.

# ExpIntVar
## The Expected Integrated Variance

- Loss function measures the overall uncertainty in the unnormalised ABC posterior over the parameter space.

- The value of the loss function depends on the next simulation so the next evaluation location $\theta*$ is chosen to minimise the expected loss

$$\theta_{t+1} = \arg \min_{\theta*} L_{1:t}(\theta*)$$

- The expected loss $L(\,\cdot\,)$ approximated as:

$$L_{1:t}(\theta*) \approx 2 \cdot \sum_{i=1}^{s} \omega^i \cdot p^2(\theta^i) \cdot w_{1:t+1}(\theta^i, \theta*)$$

- $\omega^i$ is an importance weight, $p^2(\theta^i)$ is the prior squared, and $w_{1:t+1}(\theta^i, \theta*)$ is the expected variance of the unnormalised ABC posterior at $\theta^i$ after running the simulation model with parameter $\theta*$

# Sampling from surrogate

- To represent the posterior distribution we require a set of samples drawn from it

- ABC methods produce an approximate sample from the posterior

- Surrogate methods provide an approximate posterior curve

- We use MCMC methods to draw a posterior sample

# After inference

# After inference
## How reliable are the results

- Interpret inference results

  - Different error sources

    - Algorithm performance

    - Model performance

    - Simulator performance

- Likelihood-free inference methods are based on several levels of approximations

  - All add to the total error

    - Which parts contribute most?

      - Run algorithm longer/sample more?

      - Change models?