

The School of Mathematics



THE UNIVERSITY
of EDINBURGH

Robust Optimisation Monte Carlo for Likelihood-Free Inference

by

Vasileios Gkolemis

Dissertation Presented for the Degree of
MSc in Operational Research with Data Science

August 2020

Supervised by
Dr. Michael Gutmann

Abstract

Acknowledgments

Own Work Declaration

Here comes your own work declaration

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Outline of Thesis	2
1.3	Notation	3
2	Background	5
2.1	Simulator-Based (Implicit) Models	5
2.1.1	Approximate Bayesian Computation (ABC) Rejection Sampling	5
2.1.2	Summary Statistics	5
2.1.3	Approximations Introduced	6
2.1.4	Optimization Monte Carlo (OMC)	6
2.2	Robust Optimisation Monte Carlo (ROMC) approach	7
2.2.1	Sampling and computing expectation in ROMC	8
2.2.2	Construction of the proposal region	8
2.3	Algorithmic Description of ROMC	9
2.3.1	ROMC as a Meta-Algorithm	9
2.3.2	Training and Inference Algorithms	10
2.4	Engine for Likelihood-Free Inference (ELFI) Package	11
2.4.1	Modelling	11
2.4.2	Inference Methods	12
3	Implementation	13
3.1	General Design	13
3.2	Training	13
3.3	Performing the Inference	13
3.4	Utilities	13
3.5	Computational Complexity	13
4	Experiments	14
4.1	Another Example	14
4.2	Execution Time Experiments	14
5	Conclusions	14
5.1	Outcomes	14
5.2	Future Research Directions	14
	Appendices	16
A	An Appendix	17
B	Another Appendix	18

List of Tables

List of Figures

1	Image taken from (Lintusaari et al. 2017)	2
2	Image taken from Lintusaari et al. 2018	12

laaaaaaaaaaaaaalaalaaaaakaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa

Explanation of Simulation-Based Models

Simulator-Based models are particularly captivating due to the low-level of restrictions they demand in the modeling; any physical process that can be conceptualized as a computer program of finite (deterministic or stochastic) steps, can be modelled as a Simulator-Based model without any mathematical compromise. This includes any amount of hidden (unobserved) internal variables or logic-based decisions. On the other hand, this level of freedom comes at a cost; performing inference is particularly demanding from both computational and mathematical perspective. Unfortunately, the algorithms deployed so far, allow the performance of inference only at low-dimensionality parametric spaces, i.e. $\theta \in \mathbb{R}^D$ where D is small.

Goal of Simulation-Based Models

The ROMC method (Ikonomov and Gutmann 2019) is very a recent Likelihood-Free approach; its fundamental idea is the transformation of the stochastic data generation process $M_r(\boldsymbol{\theta})$ to a deterministic mapping $g_i(\boldsymbol{\theta})$, by pre-sampling the variables that produce the randomness $\mathbf{v}_i \sim p(\mathbf{V})$. Formally, in every stochastic process the randomness is influenced by a vector of random variables \mathbf{v} , whose

1

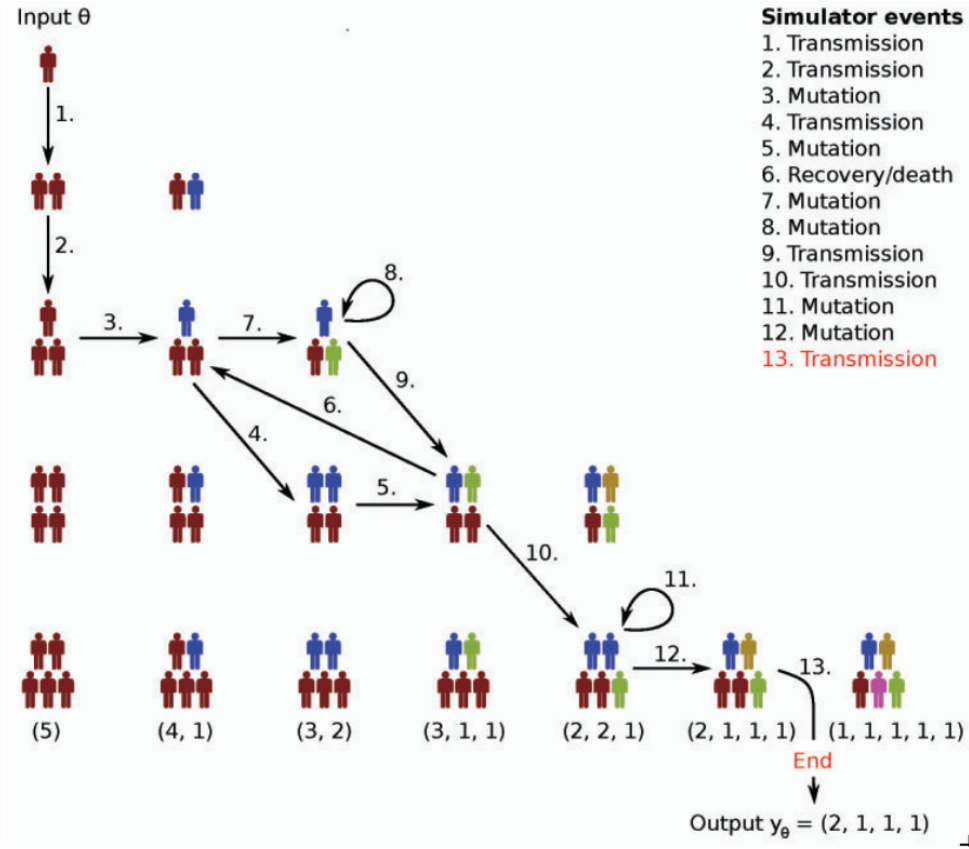


Figure 1: Image taken from (Lintusaari et al. 2017)

state is unknown before the execution of the simulation; pre-sampling this state makes the procedure deterministic, namely $g_i(\theta) = M_d(\theta, \mathbf{V} = \mathbf{v}_i)$. This approach initially introduced at (Meeds and Welling 2015) with the title *Optimisation Monte Carlo (OMC)*. The ROMC extended this approach by resolving a fundamental failure-mode of OMC. The ROMC describes a methodology for approximating the posterior through a series of steps, without explicitly enforcing which algorithms must be used for each step²; in this sense it can be thought as a meta-algorithm.

Implementation

The most important contribution of this work is the implementation of the ROMC method in the Python package Engine for Likelihood-Free Inference (ELFI) (Lintusaari et al. 2018). Since it is very recently published work the ROMC method was not implemented by now in any ML software. This works attempts to provide to the research community a tested and robust implementation for further experimentation and possible extensions.

1.2 Outline of Thesis

The remainder of the dissertation is organized as follows. In Chapter 2 we establish the mathematical formulation; more specifically we initially describe the Simulator-Based models and we provide some fundamental algorithms that have been proposed for performing statistical inference. Afterwards, we provide the mathematical description of the ROMC approach (Ikonov and Gutmann 2019). In Chapter 3, we deal with the implementation part; we initially provide some information regarding the Python package Engine for Likelihood-Free Inference (ELFI) (Lintusaari et al. 2018) and subsequently we present the implementation details of ROMC in this package. In Chapter 4, we illustrate the

²The implementation chooses a specific algorithm for each task, but this choice has just a demonstrative value; any other appropriate algorithm can be used instead.

functionalities of the ROMC implementation at some real-world examples; this chapter wishes to demonstrate the success of the ROMC method and of our implementation at Likelihood-Free tasks. Finally, in Chapter 5, we conclude with some thoughts on the work we have done and some future research ideas.

1.3 Notation

In this section, we present an overview of the symbols that will be used in the rest of the document. At this level, the quantities are introduced quite informally, through general descriptions. Most of them will be defined formally in the next chapters. We try to keep the notation as consistent as possible throughout the document. The symbol \mathbb{R}^N , when used, describes that a variable belongs to the N – *dimensional* euclidean space; N doesn't represent a specific number.

Random Generator

- $M_r(\boldsymbol{\theta}) : \mathbb{R}^D \rightarrow \mathbb{R}$: The black-box data simulator.

Parameters/Random Variables/Symbols

- $D \in \mathbb{N}$, the dimensionality of the parameter-space
- $\boldsymbol{\Theta} \in \mathbb{R}^D$, random variable representing the parameters of interest
- $\mathbf{y}_0 \in \mathbb{R}^N$, the vector of observations
- $\epsilon \in \mathbb{R}$, the threshold setting the limit on the region around \mathbf{y}_0
- $\mathbf{V} \in \mathbb{R}^N$, random variable representing the stochasticity of the generator. It is also called nuisance variable, because we are not interested in inferring a posterior distribution on it.
- $\mathbf{v}_i \sim \mathbf{V}$, a specific sample drawn from \mathbf{V}
- \mathbf{Y}_θ , random variable describing the simulator $M_r(\boldsymbol{\theta})$.
- $\mathbf{y}_i \sim \mathbf{Y}_\theta$, a sample drawn from \mathbf{Y}_θ . It can be obtained by executing the simulator $\mathbf{y}_i \sim M_r(\boldsymbol{\theta})$

Sets

- $B_{d,\epsilon}(\mathbf{y}_0)$, the set of points $\mathbf{y} := \{\mathbf{y} : d(\mathbf{y}, \mathbf{y}_0) \leq \epsilon\}$
- $B_{d,\epsilon}^i$, the set of points defined around \mathbf{y}_i i.e. $B_{d,\epsilon}^i = B_{d,\epsilon}(\mathbf{y}_i)$
- S_i , the set of $\boldsymbol{\theta}$ points, such that $\{\boldsymbol{\theta} \in B_{d,\epsilon}(\mathbf{y}_i = M_d(\boldsymbol{\theta}, \mathbf{v}_i))\}$

Generic Functions

- $p(\cdot)$, any valid pdf
- $p(\cdot|\cdot)$, any valid conditional distribution.
- $p(\boldsymbol{\theta})$, the prior distribution on the parameters
- $p(\mathbf{v})$, the prior distribution on the nuisance variables
- $p(\boldsymbol{\theta}|\mathbf{y}_0)$, the posterior distribution
- $p_{d,\epsilon}(\boldsymbol{\theta}|\mathbf{y}_0)$, the approximate posterior distribution
- $d(\mathbf{x}, \mathbf{y}) : \mathbb{R}^{2N} \rightarrow \mathbb{R}$: any valid distance, the L_2 norm: $\|\mathbf{x} - \mathbf{y}\|_2$

Functions (Mappings)

- $M_d(\boldsymbol{\theta}, \mathbf{v}) : \mathbb{R}^D \rightarrow \mathbb{R}$, the deterministic generator; all stochastic variables that are part of the data generation process are represented by the **parameter** \mathbf{v}
- $f_i(\boldsymbol{\theta}) = M_d(\boldsymbol{\theta}, \mathbf{v}_i)$, deterministic generator associated with sample $\mathbf{v}_i \sim p(\mathbf{v})$
- $g_i(\boldsymbol{\theta}) = d(f_i(\boldsymbol{\theta}), \mathbf{y}_0)$, distance of the generated data $f_i(\boldsymbol{\theta})$ from the observations
- $T(\mathbf{x}) : \mathbb{R}^{D_1} \rightarrow \mathbb{R}^{D_2}$ where $D_1 > D_2$, the mapping that computes the summary statistic
- $\mathbb{1}_{B_{d,\epsilon}(\mathbf{y}_0)}(\mathbf{y})$, the indicator function; returns 1 if $d(\mathbf{y}, \mathbf{y}_0) \leq \epsilon$, else 0
- $L(\boldsymbol{\theta})$, the likelihood
- $L_{d,\epsilon}(\boldsymbol{\theta})$, the approximate likelihood

2 Background

2.1 Simulator-Based (Implicit) Models

As already stated at Chapter 1, in simulator-based models we cannot evaluate the posterior $p(\boldsymbol{\theta}|\mathbf{y}_0) \propto L(\boldsymbol{\theta})p(\boldsymbol{\theta})$, due to the intractability of the likelihood $L(\boldsymbol{\theta}) = p(\mathbf{y}_0|\boldsymbol{\theta})$. The following equation allows incorporating the simulator in the place of the likelihood and forms the basis of all likelihood-free inference approaches.

$$L(\boldsymbol{\theta}) = \lim_{\epsilon \rightarrow 0} c_\epsilon \int_{\mathbf{y} \in B_\epsilon(\mathbf{y}_0)} p(\mathbf{y}|\boldsymbol{\theta}) d\mathbf{y} = \lim_{\epsilon \rightarrow 0} c_\epsilon \Pr(M_r(\boldsymbol{\theta}) \in B_\epsilon(\mathbf{y}_0)) \quad (2.1)$$

where c_ϵ is a proportionality factor, needed when $\Pr(M_r(\boldsymbol{\theta}) \in B_\epsilon(\mathbf{y}_0))$ tends to zero, as ϵ tends to zero. Equation 2.1 describes that the likelihood of a specific parameter configuration $\boldsymbol{\theta}$ is proportional to the probability that the simulator will produce outputs equal to the observations, given the specific configuration.

2.1.1 Approximate Bayesian Computation (ABC) Rejection Sampling

ABC rejection sampling is a modified version of the traditional rejection sampling method, for cases when the evaluation of the likelihood is intractable. In typical rejection sampling, a sample obtained from the prior $\boldsymbol{\theta} \sim p(\boldsymbol{\theta})$ gets accepted with probability $L(\boldsymbol{\theta})/\max_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$. Though we cannot use this approach out-of-the-box (evaluating $L(\boldsymbol{\theta})$ is impossible in our case), we can modify the method incorporating the simulator.

In the discrete case scenario where $\mathbf{Y}_{\boldsymbol{\theta}}$ can take a finite set of values, the likelihood becomes $L(\boldsymbol{\theta}) = \Pr(\mathbf{Y}_{\boldsymbol{\theta}} = \mathbf{y}_0)$ and the posterior $p(\boldsymbol{\theta}|\mathbf{y}_0) \propto \Pr(\mathbf{Y}_{\boldsymbol{\theta}} = \mathbf{y}_0)p(\boldsymbol{\theta})$. We can sample from the prior $\boldsymbol{\theta}_i \sim p(\boldsymbol{\theta})$, run the simulator $\mathbf{y}_i = M_r(\boldsymbol{\theta}_i)$ and accept $\boldsymbol{\theta}_i$ only if $\mathbf{y}_i = \mathbf{y}_0$.

The method above becomes less helpfull as the finite set of $\mathbf{Y}_{\boldsymbol{\theta}}$ values grows larger, since the probability of accepting a sample becomes smaller. In the limit where the set becomes infinite (i.e. continuous case) the probability becomes zero. In order for the method to work in this set-up, a relaxation is introduced; we relax the acceptance criterion by letting \mathbf{y}_i lie in a larger set of points i.e. $\mathbf{y}_i \in B_{d,\epsilon}(\mathbf{y}_0), \epsilon > 0$. The region can be defined as $B_{d,\epsilon}(\mathbf{y}_0) := \{\mathbf{y} : d(\mathbf{y}, \mathbf{y}_0) \leq \epsilon\}$ where $d(\cdot, \cdot)$ can represent any valid distance. With this modification, the maintained samples follow an approximate posterior,

$$p_{d,\epsilon}(\boldsymbol{\theta}|\mathbf{y}_0) \propto \Pr(\mathbf{y} \in B_{d,\epsilon}(\mathbf{y}_0))p(\boldsymbol{\theta}) \quad (2.2)$$

This method is called Rejection ABC.

2.1.2 Summary Statistics

When the dimensionality of $\mathbf{y} \in \mathbb{R}^D$ is high, generating samples inside $B_{d,\epsilon}(\mathbf{y}_0)$ becomes rare even when ϵ is relatively large; this is the curse of dimensionality. As a representative example lets make the following hypothesis;

- d is set to be the Euclidean distance, hence $B_{d,\epsilon}(\mathbf{y}_0) := \{\mathbf{y} : \|\mathbf{y} - \mathbf{y}_0\|_2^2 < \epsilon^2\}$ is a hyper-sphere with radius ϵ and volume $V_{hypersphere} = \frac{\pi^{D/2}}{\Gamma(D/2+1)}\epsilon^D$
- the prior $p(\boldsymbol{\theta})$ is a uniform distribution in a hyper-cube with side of length 2ϵ and volume $V_{hypercube} = (2\epsilon)^D$
- the generative model is the identity $\mathbf{y} = f(\boldsymbol{\theta}) = \boldsymbol{\theta}$

then the probability of drawing a sample inside the hypersphere equals the fraction of the volume of a hypersphere inscribed in a hypercube:

$$Pr(\mathbf{y} \in B_{d,\epsilon}(\mathbf{y}_0)) = Pr(\boldsymbol{\theta} \in B_{d,\epsilon}(\mathbf{y}_0)) = \frac{V_{hypersphere}}{V_{hypercube}} = \frac{\pi^{D/2}}{2^D \Gamma(D/2 + 1)} \rightarrow 0, \quad \text{as } D \rightarrow \infty \quad (2.3)$$

We observe that the probability tends to 0, independently of ϵ ; enlarging ϵ will not increase the acceptance rate. Intuitively, we can say that in higher-dimensions the volume of the hypercube concentrates at its corners and less volume is captured inside the hypersphere. This produces the need for a mapping $T : \mathbb{R}^{D_1} \rightarrow \mathbb{R}^{D_2}$ where $D_1 > D_2$, for squeezing the dimensionality of the output. This intermediate step redefines the area as $B_{d,\epsilon}(\mathbf{y}_0) := \{\mathbf{y} : d(T(\mathbf{y}), T(\mathbf{y}_0)) \leq \epsilon\}$. This dimensionality-reduction step is called *summary statistic* extraction, since the distance is not measured on the actual outputs, but on a summarization (i.e. lower-dimension representation) of them.

2.1.3 Approximations Introduced

So far, we have introduced some approximations for inferring the posterior as $p_{d,\epsilon}(\boldsymbol{\theta}|\mathbf{y}_0) \propto Pr(\mathbf{Y}_{\boldsymbol{\theta}} \in B_{d,\epsilon}(\mathbf{y}_0))p(\boldsymbol{\theta})$ where $B_{d,\epsilon}(\mathbf{y}_0) := \{\mathbf{y} : d(T(\mathbf{y}), T(\mathbf{y}_0)) < \epsilon\}$. These approximations introduce two different types of errors:

- ϵ is chosen to be *big enough*, so that enough samples are accepted
- T introduces loss of information, making possible a \mathbf{y} far away from the \mathbf{y}_0 i.e. $\mathbf{y} : d(\mathbf{y}, \mathbf{y}_0) > \epsilon$, to enter the acceptance region after the dimensionality reduction $d(T(\mathbf{y}), T(\mathbf{y}_0)) \leq \epsilon$

In the following sections we will not use the summary statistics in our expressions, for the notation not to clutter. One could understand it as absorbing the mapping $T(\cdot)$ as the last part of the simulator. In any case, all the following propositions are valid with the use of summary statistics.

2.1.4 Optimization Monte Carlo (OMC)

Before we define the likelihood approximation as introduced in the OMC approach, let's define the indicator function based on $B_{d,\epsilon}$.

Indicator Function

The indicator function $\mathbb{1}_{B_{d,\epsilon}(\mathbf{y})}(\mathbf{x})$ returns 1 if $\mathbf{x} \in B_{d,\epsilon}(\mathbf{y})$ and 0 otherwise. If $d(\cdot, \cdot)$ is a formal distance, due to symmetry $\mathbb{1}_{B_{d,\epsilon}(\mathbf{y})}(\mathbf{x}) = \mathbb{1}_{B_{d,\epsilon}(\mathbf{x})}(\mathbf{y})$, so they can be used interchangeably.

$$\mathbb{1}_{B_{d,\epsilon}(\mathbf{y})}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x} \in B_{d,\epsilon}(\mathbf{y}) \\ 0 & \text{otherwise} \end{cases} \quad (2.4)$$

Likelihood approximation

Based on equation (2.2) and the indicator function as defined above (2.4), we can approximate the likelihood as:

$$L_{d,\epsilon}(\boldsymbol{\theta}) = \int_{\mathbf{y} \in B_{d,\epsilon}(\mathbf{y}_0)} p(\mathbf{y}|\boldsymbol{\theta}) d\mathbf{y} = \int_{\mathbf{y} \in \mathbb{R}^D} \mathbb{1}_{B_{d,\epsilon}(\mathbf{y}_0)}(\mathbf{y}) p(\mathbf{y}|\boldsymbol{\theta}) d\mathbf{y} \quad (2.5)$$

$$\approx \frac{1}{N} \sum_i^N \mathbb{1}_{B_{d,\epsilon}(\mathbf{y}_0)}(\mathbf{y}_i), \text{ where } \mathbf{y}_i \sim M_r(\boldsymbol{\theta}) \quad (2.6)$$

$$\approx \frac{1}{N} \sum_i^N \mathbb{1}_{B_{d,\epsilon}(\mathbf{y}_0)}(\mathbf{y}_i) \text{ where } \mathbf{y}_i = M_d(\boldsymbol{\theta}, \mathbf{v}_i), \mathbf{v}_i \sim p(\mathbf{v}) \quad (2.7)$$

This approach is quite intuitive; approximating the likelihood of a specific θ requires sampling from the data generator and count the fraction of samples that lie inside the area around the observations. Nevertheless, using the approximation of equation (2.6) we need to draw N new samples for each distinct evaluation of $L_{d,\epsilon}(\theta)$; this makes this approach quite inconvenient from a computational point-of-view. For this reason, we choose to approximate the integral as in equation (2.7); the nuisance variables are sampled once $\mathbf{v}_i \sim p(\mathbf{v})$ and we count the fraction of samples that lie inside the area using the deterministic simulator $M_d(\theta, \mathbf{v}_i)$. Hence, the evaluation of each θ does not imply drawing any new samples. Based on this approach, the unnormalized approximate posterior can be defined as:

$$p_{d,\epsilon}(\theta|\mathbf{y}_0) \propto p(\theta) \sum_i^N \mathbb{1}_{B_{d,\epsilon}(\mathbf{y}_0)}(\mathbf{y}_i) \quad (2.8)$$

Further approximations for sampling and computing expectations

The posterior approximation in (2.2) does not provide any obvious way for drawing samples. In fact, the set $S_i = \{\theta \in B_{d,\epsilon}(\mathbf{y}_i = M_d(\theta, \mathbf{v}_i))\}$ can represent any arbitrary shape in the D -dimensional Euclidean space; it can be non-convex, can contain disjoint sets of θ etc. This observation leads to the need for a further simplification of the posterior for sampling from it.

As a side-note, sampling could be performed in a straightforward fashion using importance sampling; using the prior as the proposal distribution $\theta_i \sim p(\theta)$ and computing the weight as $w_i = \frac{L_{d,\epsilon}(\theta_i)}{p(\theta_i)}$. This approach has the same drawbacks as ABC rejection sampling; when the prior is wide or the dimensionality D is high, drawing a sample with non-zero weight is rare, leading to either poor Effective Sample Size (ESS) or huge execution time.

The OMC proposes a quite drastic simplification of the posterior; it squeezes all regions S_i into a single point $\theta_i^* \in S_i$ attaching a weight w_i proportional to the volume of S_i . For obtaining θ_i^* , a gradient based optimiser is used for minimising $g_i(\theta) = d(\mathbf{y}_0, f_i(\theta))$ and the estimation of the volume of S_i is done using the Hessian approximation $\mathbf{H}_i \approx \mathbf{J}_i^{*T} \mathbf{J}_i^*$. Hence,

$$p(\theta|\mathbf{y}_0) \propto p(\theta) \sum_i^N w_i \delta(\theta - \theta_i^*) \quad (2.9)$$

$$\theta_i^* = \operatorname{argmin}_{\theta} [g_i(\theta)] \quad (2.10)$$

$$w_i \propto \frac{1}{\sqrt{\det(\mathbf{J}_i^{*T} \mathbf{J}_i^*)}} \quad (2.11)$$

The distribution (2.9) provides automatically weighted samples and an expectation can be computed easily with the following equation,

$$E_{p(\theta|\mathbf{y}_0)}[h(\theta)] = \frac{\sum_i^N w_i p(\theta_i^*) h(\theta_i^*)}{\sum_i^N w_i p(\theta_i^*)} \quad (2.12)$$

2.2 Robust Optimisation Monte Carlo (ROMC) approach

The simplifications introduced by OMC, although quite usefull from a computational point of view, they suffer from some major failure modes:

- The whole set of acceptable S_i for each nuisance variable gets shrinked to a single point θ_i^* ; this can add significant error when then the area S_i is relatively big.
- The weight w_i is computed using only **local** information, i.e. the curvature of g_i at the point θ_i^* . This approach can introduce significant error when g_i is almost flat at θ_i^* , leading to a $\det(\mathbf{J}_i^{*T} \mathbf{J}_i^*) \rightarrow 0 \Rightarrow w_i \rightarrow \infty$, dominating the posterior.

- There is no way to solve the optimisation problem $\theta_i^* = \operatorname{argmin}_{\theta} [g_i(\theta)]$ when g_i is not differentiable.

2.2.1 Sampling and computing expectation in ROMC

The ROMC approach resolves the aforementioned issues. Instead of collapsing the acceptance regions into single points, it tries to approximate them with a bounding box and then define a uniform distribution on them.³ This distribution serves as a proposal distribution for importance sampling. If we define as q_i , the uniform distribution defined on the i -th bounding box, weighted sampling is performed as:

$$\theta_{ij} \sim q_i \quad (2.13)$$

$$w_{ij} = \frac{\mathbb{1}_{B_{d,\epsilon}(\mathbf{y}_0)}(M_d(\theta_{ij}, \mathbf{v}_i))p(\theta_{ij})}{q(\theta_{ij})} \quad (2.14)$$

Having defined the procedure for obtaining weighted samples, any expectation $E_{p(\theta|\mathbf{y}_0)}[h(\theta)]$, can be approximated as,

$$E_{p(\theta|\mathbf{y}_0)}[h(\theta)] \approx \frac{\sum_{ij} w_{ij} h(\theta_{ij})}{\sum_{ij} w_{ij}} \quad (2.15)$$

2.2.2 Construction of the proposal region

In this section we will describe mathematically the steps needed for computing the proposal distributions q_i . There will be also presented a Bayesian Optimisation alternative when gradients are not available.

Define and solve deterministic optimisation problems

For each set of nuisance variables $\mathbf{v}_i, i = \{1, 2, \dots, n_1\}$ a deterministic function is defined as $f_i(\theta) = M_d(\theta, \mathbf{v}_i)$. For constructing the proposal region, we search for a point $\theta_* : d(f_i(\theta_*), \mathbf{y}_0) < \epsilon$; this point can be obtained by solving the the following optimisation problem:

$$\min_{\theta} \quad g_i(\theta) = d(\mathbf{y}_0, f_i(\theta)) \quad (2.16a)$$

$$\text{subject to} \quad g_i(\theta) \leq \epsilon \quad (2.16b)$$

We maintain a list of the solutions θ_i^* of the optimisation problems. If for a specific set of nuisance variables \mathbf{v}_i , there is no feasible solution we add nothing to the list. The optimisation problem can be treated as unconstrained, accepting the optimal point $\theta_i^* = \operatorname{argmin}_{\theta} g_i(\theta)$ only if $g_i(\theta_i^*) < \epsilon$.

Gradient-Based Approach

The nature of the generative model $M_r(\theta)$, specifies the properties of the objective function g_i . If g_i is continuous with smooth gradients $\nabla_{\theta} g_i$ any gradient-based iterative algorithm can be used for solving 2.16a. The gradients $\nabla_{\theta} g_i$ can be either provided in closed form or approximated by finite differences.

Bayesian Optimisation Approach

In cases where the gradients are not available, the Bayesian Optimisation scheme provides an alternative choice (Shahriari et al. 2016). With this approach, apart from obtaining an optimal θ_i^* , a surrogate model \hat{d}_i of the distance g_i is fitted; this approximate model can be used in the following

³The description on how to estimate the Bounding Box is given in the following chapters.

Algorithm 1 ROMC as a Meta-Algorithm. Requires $M_d(\theta, \mathbf{v}), \mathbf{y}_0$. Hyperparameters n_1, n_2 .

```

1: for  $i \leftarrow 1$  to  $n_1$  do
2:   Sample a random state  $\mathbf{v}_i \sim p(\mathbf{v})$ 
3:   Define the deterministic mapping  $f_i(\theta) = M_d(\theta, \mathbf{v})$  and therefore  $g_i(\theta) = d(f_i(\theta), \mathbf{y}_0)$ .
4:   Obtain  $d_i^* = \min_{\theta} g_i(\theta)$  and  $\theta_i^* = \operatorname{argmin}_{\theta} g_i(\theta)$  using any convenient optimiser.
5:   Approximate the local area  $\{\theta : g_i(\theta) < \epsilon \text{ and } d(\theta, \theta_i^*) < M\}$  with a bounding box, using any
   convenient method.
6:   Define a proposal distribution  $q_i(\theta)$ , exploiting the Bounding Box.
7:   for  $j \leftarrow 1$  to  $n_2$  do
8:      $\theta_{ij} \sim q_i(\theta)$ 
9:     Accept  $\theta_{ij}$  as posterior sample with weight  $w_{ij} = \frac{p(\theta_{ij})}{q_i(\theta_{ij})} \mathbb{1}_{B_{d,\epsilon}^i}(\theta_{ij})$ 
return (List with samples  $\theta_{ij}$  and weights  $w_{ij}$ )

```

steps for making the method more efficient. Specifically, in the construction of the proposal region and in equations (2.2), (2.13), (2.15) it could replace g_i in the evaluation of the indicator function 2.4, providing a major speed-up.

Construction of the proposal area q_i

After obtaining θ_i^* such that $g_i(\theta_i^*) < \epsilon$, we need to construct a bounding box around it. The bounding box must contain the acceptance region around θ_i^* , i.e. $\{\theta : g_i(\theta) < \epsilon \text{ and } d(\theta, \theta_i^*) < M\}$. The second condition $d(\theta, \theta_i^*) < M$ is meant to describe that if $S_i := \{\theta : g_i(\theta) < \epsilon\}$ contains a number of disjoint sets of θ that respect $g_i(\theta) < \epsilon$, we want our bounding box to fit the one that contains θ_i^* . We seek for a bounding box to be as tight as possible to the local acceptance region (enlarging the bounding box decreases the acceptance rate) but large enough for not discarding accepted areas.

In contrast with the OMC approach, we construct the bounding box by querying the indication function along the search directions. The search directions \mathbf{v}_d are computed as the eigenvectors of the curvature at θ_i^* and a line-search method is used to obtain the limit point where $g_i(\theta_i^* + \kappa \mathbf{v}_d) \geq \epsilon$. Algorithm 4 describes analytically the method. After the limits are obtained along all search directions, we define the uniform distribution q_i on the bounding box. This is the proposal distribution of the importance sampling explained in (2.13).

2.3 Algorithmic Description of ROMC

In this section, we attempt the depiction of the mathematical description of ROMC in an algorithmic view. Specifically, in section 2.3.1 we present the general algorithmic description of ROMC as a meta-algorithm and in section 2.3.2 the proposals of ROMC for solving the training and the inference parts.

2.3.1 ROMC as a Meta-Algorithm

As stated at the introductory chapter, ROMC can be understood as step-by-step algorithmic approach for performing the inference in Simulator-Based Models. The particular methods used for solving the sub-tasks are left as a free choice to the user. As presented in Algorithm 1, the methods involved in solving the optimisation problem (step 4) and constructing the bounding box (step 5) are not restricted. The practitioner may choose any convenient algorithm, judging the trade-offs between accuracy, robustness, efficiency and complexity. In particular for the optimisation step, the choice of the appropriate optimiser should also consider the properties of $g_i(\theta)$. Some important questions that should be considered are whether the function differentiable and if so whether we know the gradients $\nabla_{\theta} g_i$ in closed-form. As described in sections 2.2.2 and 2.2.2, the ROMC paper considers two alternative optimisation schemes (gradient-based and Bayesian-optimisation approach) depending on whether the gradients are available or not.

2.3.2 Training and Inference Algorithms

In this section, we will provide the algorithmic description of the ROMC method; (a) the procedures for solving the optimisation problems using either the gradient based approach or the Gaussian Process alternative and (b) the construction of the Bounding Box. Afterwards, we will discuss the advantages and the disadvantages of each choice both in terms of accuracy and efficiency. At a high-level, the ROMC method can be split into the training and the inference part.

At the training (fitting) part, the goal is the estimation of the proposal regions q_i . The steps include (a) sampling the nuisance variables $\mathbf{v}_i \sim p(\mathbf{v})$ (b) defining the optimisation problems $\min_{\boldsymbol{\theta}}[g_i(\boldsymbol{\theta})]$ (c) obtaining $\boldsymbol{\theta}_i^*$ (d) checking whether $d_i^* \leq \epsilon$ and (e) building the bounding box for obtaining the proposal region q_i . If gradients are available, using a gradient-based method is advised for obtaining $\boldsymbol{\theta}_i^*$ much faster. Providing $\nabla_{\boldsymbol{\theta}} g_i$ in closed-form provides an upgrade in both accuracy and efficiency; If closed-form description is not available, approximate gradients with finite-differences requires two evaluations of g_i for **every** parameter $\boldsymbol{\theta}_d$. For low-dimensional problems though, this approach still works well. When gradients are not available or g_i is not differentiable, using the Bayesian Optimisation Paradigm is a solution. In this case, the training part is much slower due to the fitting of the surrogate model and the ignorance of the slope throughout the optimisation procedure. Nevertheless, computing the proposal region q_i becomes faster since \hat{d}_i can be used instead of g_i which involves running the whole simulator $M_d(\boldsymbol{\theta}, \mathbf{v}_i)$ for each query. The algorithms 2 and 3 present the above procedure.

Performing the inference includes one or more of the following three tasks; (a) evaluating the unnormalised posterior $p_{d,\epsilon}(\theta_b|\mathbf{y}_0)$ (b) sampling from the posterior $\boldsymbol{\theta}_i \sim p_{d,\epsilon}(\theta_b|\mathbf{y}_0)$ (c) computing an expectation $E_{\boldsymbol{\theta}|\mathbf{y}_0}[h(\boldsymbol{\theta})]$. Computing an expectation can be done easily after weighted samples are obtained using the equation 2.15, so we will not discuss it separately.

Evaluating the unnormalised posterior requires solely the deterministic functions g_i and the prior distribution $p(\boldsymbol{\theta})$; there is no need for solving the optimisation problems and building the proposal regions. The evaluation requires iterating over all g_i and evaluating the distance from the observed data. In contrast, using the GP approach, the optimisation part should be performed first for fitting the surrogate models $\hat{d}_i(\boldsymbol{\theta})$ and evaluate the indicator function on them. This provides an important speed-up, especially when running the simulator is computationally expensive.

Sampling is performed by getting n_2 samples from each proposal distribution q_i . For each sample $\boldsymbol{\theta}_{ij}$, the indicator function is evaluated $\mathbb{1}_{B_{d,\epsilon}^i(\mathbf{y}_0)}(\boldsymbol{\theta}_{ij})$ for checking if it lies inside the acceptance region. If so the corresponding weight is computed as in (2.13). As before, if a surrogate model \hat{d} has been fitted, it can be used for the evaluation of the indicator function providing again a speedup. Apparently, the computational benefit is more important compared to posterior evaluation, because the indicator function must be evaluated for a total of $n_1 \times n_2$ points. The sampling algorithms are presented step-by-step in algorithms 5 and 6.

In summary, we can state that the choice of using a bayesian optimisation approach provides a significant speed-up in the inference part with the cost of making the training part slower and a possible approximation error. It is typical in many Machine-Learning use cases, being able to provide enough time and computational resources for the training phase, but asking for efficiency in the inference part.

Algorithm 2 Training Part - Gradient approach. Requires $g_i(\boldsymbol{\theta}), p(\boldsymbol{\theta})$

```

1: for  $i \leftarrow 1$  to  $n$  do
2:   Obtain  $\boldsymbol{\theta}_i^*$  using a Gradient Optimiser
3:   if  $g_i(\boldsymbol{\theta}_i^*) > \epsilon$  then
4:     go to 1
5:   else
6:     Approximate  $\mathbf{J}_i^* = \nabla g_i(\boldsymbol{\theta})$  and  $H_i \approx \mathbf{J}_i^T \mathbf{J}_i$ 
7:   Use Algorithm 4 to obtain  $q_i$ 
return  $q_i, p(\boldsymbol{\theta}), g_i(\boldsymbol{\theta})$ 
```

Algorithm 3 Training Part - GP approach. Requires $g_i(\boldsymbol{\theta}), p(\boldsymbol{\theta})$

```

1: for  $i \leftarrow 1$  to  $n$  do
2:   Obtain  $\boldsymbol{\theta}_i^*, \hat{d}_i(\boldsymbol{\theta})$  using a GP approach
3:   if  $g_i(\boldsymbol{\theta}_i^*) > \epsilon$  then
4:     go to 1
5:   else
6:     Approximate  $H_i \approx \mathbf{J}_i^T \mathbf{J}_i$ 
7:   Use Algorithm 4 to obtain  $q_i$ 
return  $q_i, p(\boldsymbol{\theta}), \hat{d}_i(\boldsymbol{\theta})$ 
```

Algorithm 4 Computation of the proposal distribution q_i ; Needs, a model of distance d , optimal point θ_i^* , number of refinements K , step size η and curvature matrix \mathbf{H}_i ($\mathbf{J}_i^T \mathbf{J}_i$ or GP Hessian)

```

1: Compute eigenvectors  $\mathbf{v}_d$  of  $\mathbf{H}_i$  ( $d = 1, \dots, ||\theta||$ )
2: for  $d \leftarrow 1$  to  $||\theta||$  do
3:    $\tilde{\theta} \leftarrow \theta_i^*$ 
4:    $k \leftarrow 0$ 
5:   repeat
6:     repeat
7:        $\tilde{\theta} \leftarrow \tilde{\theta} + \eta \mathbf{v}_d$  ▷ Large step size  $\eta$ .
8:     until  $d(f_i(\tilde{\theta}), \mathbf{y}_0) > \epsilon$ 
9:      $\tilde{\theta} \leftarrow \tilde{\theta} - \eta \mathbf{v}_d$ 
10:     $\eta \leftarrow \eta/2$  ▷ More accurate region boundary
11:     $k \leftarrow k + 1$ 
12:  until  $k = K$ 
13:  Set final  $\tilde{\theta}$  as region end point.
14:  Repeat steps 3 - 13 for  $\mathbf{v}_d = -\mathbf{v}_d$ 
15: Fit a rectangular box around the region end points and define  $q_i$  as uniform distribution

```

Algorithm 5 Sampling - Gradient Based approach. Requires $g_i(\theta), p(\theta), q_i$

```

1: for  $i \leftarrow 1$  to  $n_1$  do
2:   for  $j \leftarrow 1$  to  $n_2$  do
3:      $\theta_{ij} \sim q_i$ 
4:     if  $g_i(\theta_{ij}) > \epsilon$  then
5:       Reject  $\theta_{ij}$ 
6:     else
7:        $w_{ij} = \frac{p(\theta_{ij})}{q(\theta_{ij})}$ 
8:       Accept  $\theta_{ij}$ , with weight  $w_{ij}$ 

```

Algorithm 6 Sampling - GP approach. Requires $\hat{d}_i(\theta), p(\theta), q_i$

```

1: for  $i \leftarrow 1$  to  $n_1$  do
2:   for  $j \leftarrow 1$  to  $n_2$  do
3:      $\theta_{ij} \sim q_i$ 
4:     if  $\hat{d}_i(\theta_{ij}) > \epsilon$  then
5:       Reject  $\theta_{ij}$ 
6:     else
7:        $w_{ij} = \frac{p(\theta_{ij})}{q(\theta_{ij})}$ 
8:       Accept  $\theta_{ij}$ , with weight  $w_{ij}$ 

```

2.4 Engine for Likelihood-Free Inference (ELFI) Package

!! NOT fully written, I have to add some more info !!

The Engine for Likelihood-Free Inference (ELFI) Lintusaari et al. 2018 is a Python software library dedicated to likelihood-free inference (LFI). ELFI models in a convenient manner all the fundamental components of a Probabilistic Model such as priors, simulators, summaries and distances. Furthermore, in the ELFI there are implemented a wide range of likelihood-free inference methods.

2.4.1 Modelling

ELFI models the Probabilistic Model as Directed Acyclic Graph (DAG); it implements this functionality based on the package NetworkX, which is designed for creating general purpose graphs. Although not restricted to that, in most cases the structure of a likelihood-free model follows the pattern presented in figure 2; there are edges that connect the *prior* distributions to the simulator, the simulator is connected to the summary statistics which are connected to the distance, which is the output node. Samples can be obtained from all nodes through sequential sampling. The nodes that are defined as *elfi.Prior*⁴ are automatically considered as the parameters of interest and they are the only nodes that should provide pdf evaluation, apart from sampling. The function passed as argument in the *elfi.Summary* node can be any valid Python function with arguments the prior variables.

⁴The *elfi.Prior* functionality is a wrapper around the *scipy.stats* package.

```

# Define the simulator, the summary and the observed data
def simulator(t1, t2, batch_size=1, random_state=None):
    # Implementation comes here. Return 'batch_size'
    # simulations wrapped to a NumPy array.
def summary(data, argument=0):
    # Implementation comes here...
y = # Observed data, as one element of a batch.

# Specify the ELFI graph
t1 = elfi.Prior('uniform', -2, 4)
t2 = elfi.Prior('normal', t1, 5) # depends on t1
SIM = elfi.Simulator(simulator, t1, t2, observed=y)
S1 = elfi.Summary(summary, SIM)
S2 = elfi.Summary(summary, SIM, 2)
d = elfi.Distance('euclidean', S1, S2)

# Run the rejection sampler
rej = elfi.Rejection(d, batch_size=10000)
result = rej.sample(1000, threshold=0.1)

```

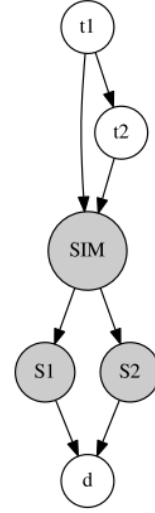


Figure 2: Image taken from Lintusaari et al. 2018

2.4.2 Inference Methods

All Inference Methods that are implemented in ELFI, follow some common guidelines; (a) their initialization should be defined by passing the output graph as the initial argument and afterwards come the rest hyperparameters of the method and (b) they must provide a basic inference functionality, e.g. `<method>.sample()`, which returns a predefined `elfi.Result` object containing the obtained samples along with some other useful functionalities (e.g. plotting the marginal posteriors).

A good collection of likelihood-free inference methods is implemented so far, such as the *ABC Rejection Sampler* and *Sequential Monte Carlo ABC Sampler*. A quite central method implemented by ELFI is the *Bayesian Optimization for Likelihood-Free Inference (BOLFI)*, which is methodologically quite close to the ROMC method we implement in the current dissertation.

3 Implementation

3.1 General Design

3.2 Training

3.3 Performing the Inference

3.4 Utilities

3.5 Computational Complexity

4 Experiments

4.1 Another Example

4.2 Execution Time Experiments

5 Conclusions

5.1 Outcomes

5.2 Future Research Directions

References

- [CG19] Yanzhi Chen and Michael U Gutmann. “Adaptive Gaussian Copula ABC”. In: *Proceedings of Machine Learning Research*. Vol. 89. 2019, pp. 1584–1592. URL: <http://proceedings.mlr.press/v89/chen19d.html>.
- [GC16] Michael U. Gutmann and Jukka Corander. *Bayesian optimization for likelihood-free inference of simulator-based statistical models*. 2016. arXiv: 1501.03291.
- [IG19] Borislav Ikononov and Michael U. Gutmann. “Robust Optimisation Monte Carlo”. In: (2019). arXiv: 1904.00670. URL: <http://arxiv.org/abs/1904.00670>.
- [Lin+17] Jarno Lintusaari et al. “Fundamentals and recent developments in approximate Bayesian computation”. In: *Systematic Biology* 66.1 (2017), e66–e82. ISSN: 1076836X. DOI: 10.1093/sysbio/syw077.
- [Lin+18] Jarno Lintusaari et al. *ELFI: Engine for Likelihood Free Inference*. 2018. eprint: arXiv: 1708.00707.
- [MW15] Edward Meeds and Max Welling. “Optimization Monte Carlo: Efficient and embarrassingly parallel likelihood-free inference”. In: *Advances in Neural Information Processing Systems*. 2015. arXiv: 1506.03693.
- [Sha+16] Bobak Shahriari et al. *Taking the human out of the loop: A review of Bayesian optimization*. 2016. DOI: 10.1109/JPROC.2015.2494218.
- [Tan+06] Mark M. Tanaka et al. “Using approximate bayesian computation to estimate tuberculosis transmission parameters from genotype data”. In: *Genetics* (2006). ISSN: 00166731. DOI: 10.1534/genetics.106.055574.

Appendices

A An Appendix

Some stuff.

B Another Appendix

Some other stuff.