

# THE UNIVERSITY of EDINBURGH

## Edinburgh Research Explorer

### **Robust Optimisation Monte Carlo**

#### Citation for published version:

Ikonomov, B & Gutmann, M 2020, Robust Optimisation Monte Carlo. in Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics (AISTATS) 2020. Proceedings of Machine Learning Research, vol. 108, PMLR, pp. 2819-2829, 23rd International Conference on Artificial Intelligence and Statistics, Online, Italy, 26/08/20. <a href="http://proceedings.mlr.press/v108/ikonomov20a.html">http://proceedings.mlr.press/v108/ikonomov20a.html</a>

Link: Link to publication record in Edinburgh Research Explorer

**Document Version:** Peer reviewed version

**Published In:** Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics (AISTATS) 2020

#### **General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



#### **Robust Optimisation Monte Carlo**

#### **Borislav Ikonomov**

School of Informatics University of Edinburgh borislav.ikonomov@ed.ac.uk

#### Abstract

This paper is on Bayesian inference for parametric statistical models that are defined by a stochastic simulator which specifies how data is generated. Exact sampling is then possible but evaluating the likelihood function is typically prohibitively expensive. Approximate Bayesian Computation (ABC) is a framework to perform approximate inference in such situations. While basic ABC algorithms are widely applicable, they are notoriously slow and much research has focused on increasing their efficiency. Optimisation Monte Carlo (OMC) has recently been proposed as an efficient and embarrassingly parallel method that leverages optimisation to accelerate the inference. In this paper, we demonstrate an important previously unrecognised failure mode of OMC: It generates strongly overconfident approximations by collapsing regions of similar or near-constant likelihood into a single point. We propose an efficient, robust generalisation of OMC that corrects this. It makes fewer assumptions, retains the main benefits of OMC, and can be performed either as post-processing to OMC or as a stand-alone computation. We demonstrate the effectiveness of the proposed Robust OMC on toy examples and tasks in inverse-graphics where we perform Bayesian inference with a complex image renderer.

#### 1 Introduction

Simulator-based models can describe many complex processes that occur in nature, such as the evolution of genomes (Marttinen et al., 2015) or the dynamics of gene regulation (Toni et al., 2009). Learning their parameters, in particular when done in a Bayesian framework, allows us to make predictions or take decisions based on incom-

#### Michael U. Gutmann

School of Informatics University of Edinburgh michael.gutmann@ed.ac.uk

plete information. However, learning the parameters or obtaining their posterior distribution is typically computationally very demanding as their likelihood functions are intractable. Likelihood-Free Inference (LFI) methods have thus emerged that perform inference when the likelihood function is not available in closed form but sampling from the model is possible.

A prominent instance of LFI is Approximate Bayesian Computation (ABC); for recent reviews, see for example (Sisson et al., 2018; Lintusaari et al., 2017). Other instances of LFI are the synthetic likelihood approach by Wood (2010) and its generalisations (Thomas et al., 2016; Price et al., 2017; Fasiolo et al., 2018). This paper focuses on ABC where the basic idea is to identify the parameter values which generate synthetic data that is close to the observed data under some chosen discrepancy measure. This measure can be the Euclidean distance between suitably chosen summary statistics, but other measures are possible too (e.g. Gutmann et al., 2014; Bernton et al., 2018). Generally, there are two main avenues of research for ABC — one focuses on improving the distance metric and/or the summary statistics used (e.g. Aeschbacher et al., 2012; Fearnhead and Prangle, 2012), while the other concentrates on computational efficiency (e.g. Beaumont et al., 2002; Blum et al., 2013; Meeds and Welling, 2015; Gutmann and Corander, 2016; Papamakarios and Murray, 2016). This paper focuses on the latter, and as such we assume the distance and summary statistics are given.

The primary focus of this paper is Optimisation Monte Carlo (OMC) — an ABC method developed by Meeds and Welling (2015) and also independently by Forneron and Ng (2016, 2018) under the name of "the reverse sampler". It uses optimisation to efficiently produce weighted posterior samples in a fully parallelisable manner, which makes it a desirable ABC method.

A weight produced by OMC represents the volume of the parameter region around a posterior sample which contains points that are as good as the sample. These points should thus be considered to be samples from the posterior too. However, if this region is particularly big, it is no longer appropriate to approximate the entire region with a single point and, as a result, OMC produces an overly confident posterior. Figure 1 illustrates this failure case for a simple 1D scenario. We can see that OMC fails to characterise the posterior uncertainty and collapses regions of similar likelihood into a single point.



Figure 1: An example where OMC fails to approximate the true posterior, collapsing a region of similar likelihood into a single point. Heuristic OMC is a simple heuristic that (unsuccessfully) attempts to solve this issue. Robust OMC is the approach proposed in this paper. See Subsection 4.1 for details.

We propose Robust OMC (ROMC) which efficiently identifies the regions of acceptable points themselves and samples from them directly. These regions are characterised by the set of all parameters which generate data whose distance to the observed data is less than a certain threshold. Instead of approximating a possibly very large region with a single point, we draw samples from the region which then replace the original OMC sample along with its weight. The main improvements of ROMC over OMC are:

- It handles likelihoods that are (nearly) flat on significant regions in parameter space. All our experiments confirm this.
- It can be applied as post-processing to an original OMC run (see experiment 2), or as part of standalone run (see experiment 3).
- OMC requires that the derivatives of the simulator can be computed or reasonably approximated, while ROMC does not (see experiment 3).

#### 2 Background

We here present the basics of Approximate Bayesian Computation (ABC), review the OMC algorithm by Meeds and Welling (2015), and discuss when it collapses regions of similar likelihood into a single point.

#### 2.1 Rejection ABC

ABC methods produce samples from an approximate posterior (see e.g. Lintusaari et al., 2017). They generally make minimal assumptions about the model and only assume black-box access to the simulator  $g(\theta, \mathbf{u})$ . The target parameters  $\theta$ , which we wish to infer, are used as input to the simulator which then stochastically produces synthetic data by drawing  $\mathbf{u}$  using a random number generator. These  $\mathbf{u}$  correspond to nuisance variables since we do not aim to find a distribution over them.

The simplest ABC method is Rejection ABC. In each iteration, data  $\mathbf{x}_i$  is simulated using the generative model  $g(\boldsymbol{\theta}, \mathbf{u}_i)$  for some setting  $\mathbf{u}_i$  of the nuisance variables and with parameter values sampled from the prior. The distance between the simulated and observed data  $d_i = d(\mathbf{x}_i, \mathbf{x}^0)$  is then computed and stored. After a sufficiently large amount of samples has been generated, the algorithm accepts those having the *n* lowest  $d_i$  as samples from the approximate posterior. See Algorithm 4 in the appendix for pseudo-code.

While simple and robust, Rejection ABC is known to be computationally inefficient, especially when the prior space is large (e.g. Lintusaari et al., 2017). In most cases, more sophisticated methods are necessary.

#### 2.2 Optimisation Monte Carlo

We start our brief review of Optimisation Monte Carlo (OMC) by noting that ABC algorithms in general implicitly approximate the likelihood function by the probability  $\Pr(d(g(\theta, \mathbf{u}), \mathbf{x}^0) \leq \epsilon)$  that the generated data is within distance  $\epsilon$  of the observed data (e.g. Lintusaari et al., 2017). ABC algorithms thus produce samples from the following approximate posterior

$$p_{\epsilon}^{*}(\boldsymbol{\theta}|\mathbf{x}^{0}) \propto p(\boldsymbol{\theta}) \operatorname{Pr}(d(g(\boldsymbol{\theta},\mathbf{u}),\mathbf{x}^{0}) \leq \epsilon |\boldsymbol{\theta}) \quad (1)$$

$$\propto p(\boldsymbol{\theta}) \int p(\mathbf{u}) \mathbb{1}_{C_{\epsilon}}(\boldsymbol{\theta}, \mathbf{u}) \, d\mathbf{u},$$
 (2)

where  $p(\mathbf{u})$  is the density of the nuisance variables  $\mathbf{u}$ ,  $C_{\epsilon} = \{(\boldsymbol{\theta}, \mathbf{u}) : d(g(\boldsymbol{\theta}, \mathbf{u}), \mathbf{x}^0) \leq \epsilon\}$  is the set of points  $(\boldsymbol{\theta}, \mathbf{u})$  for which the distance is below the threshold, and  $\mathbb{1}_{C_{\epsilon}}(\boldsymbol{\theta}, \mathbf{u})$  is an indicator function that equals one only if  $(\boldsymbol{\theta}, \mathbf{u}) \in C_{\epsilon}$ . While this formulation uses the indicator function (boxcar kernel), more general kernels can be used as well.

The integral over **u** corresponds to an expectation with respect to  $p(\mathbf{u})$  and can thus be approximated as a sample average so that we obtain the approximation

$$p_{\epsilon}(\boldsymbol{\theta}|\mathbf{x}^{0}) \propto p(\boldsymbol{\theta}) \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{C_{\epsilon}}(\boldsymbol{\theta}, \mathbf{u}_{i}),$$
 (3)

where the  $\mathbf{u}_i$  are sampled from  $p(\mathbf{u})$ . Importantly, this formulation essentially removes the randomness from the simulator:  $g(\boldsymbol{\theta}, \mathbf{u}_i)$ , which occurs in each  $\mathbb{1}_{C_{\epsilon}}(\boldsymbol{\theta}, \mathbf{u}_i)$ by definition of  $C_{\epsilon}$ , is a deterministic function of  $\boldsymbol{\theta}$  because  $\mathbf{u}_i$  is held fixed.

OMC exploits the fact that  $g(\theta, \mathbf{u}_i)$  is a deterministic function in order to accelerate the sampling from the posterior in (3). For each  $\mathbf{u}_i$ , OMC finds a value  $\theta_i^*$  for which  $g(\theta_i^*, \mathbf{u}_i)$  and  $\mathbf{x}^0$  are within distance  $\epsilon$ . Importantly, this is done by minimising the deterministic cost function  $d(g(\theta, \mathbf{u}_i), \mathbf{x}^0)$  with respect to  $\theta$ . Note that settings of  $\mathbf{u}_i$  for which no  $\theta_i^*$  has distance below  $\epsilon$  are excluded from the sum and thus do not affect the posterior.

Meeds and Welling (2015) consider the case where the distance  $d(g(\theta, \mathbf{u}), \mathbf{x}^0))$  is the Euclidean distance between some summary statistics  $\Phi$  of the generated and observed data. We denote the summary statistics of the observed data by  $\mathbf{y}^0 = \Phi(\mathbf{x}^0)$  and we further absorb the computation of the summary statistics into the simulator so that we obtain  $\mathbf{f}(\theta, \mathbf{u}) = \Phi(g(\theta, \mathbf{u}))$ , which can be regarded as a generative model on the level of the summary statistics. With this notation, OMC considers the distance  $d(g(\theta, \mathbf{u}), \mathbf{x}^0) = ||\mathbf{f}(\theta, \mathbf{u}) - \mathbf{y}^0||$ .

OMC approximates the posterior  $p_{\epsilon}^*$  in the limit of  $\epsilon \to 0$ as a mixture of weighted point masses centred at the minimisers  $\theta_i^*$ :  $p(\theta|\mathbf{x}^0) \propto \sum_{i=1}^n w_i \delta(\theta - \theta_i^*)$ . The value  $w_i$  is a weight that reflects the local behaviour of the distance function and hence  $\mathbf{f}(\theta, \mathbf{u}_i)$  around  $\theta_i^*$ . As shown by Meeds and Welling (2015), it equals  $p(\theta_i^*) *$  $\det(\mathbf{J}_i^{\mathsf{T}}\mathbf{J}_i)^{-1/2}$  where  $\mathbf{J}_i$  is the Jacobian matrix with columns  $\partial \mathbf{f}(\theta_i^*, \mathbf{u}_i)/\partial \theta_k$  with  $\theta_k$  denoting the k-th element of  $\theta$ . Algorithm 1 summarises the OMC algorithm. For further details, we refer the reader to the original paper by Meeds and Welling (2015).

Algorithm 1 Optimisation Monte Carlo. Generates n independent samples  $\theta_i^*$  with weights  $w_i$  from the approximate posterior.

1: for  $i \leftarrow 1$  to n do 2:  $\mathbf{u}_i \sim p(\mathbf{u}) \Rightarrow \text{Draw nuisance parameters } \mathbf{u}_i$ . 3:  $\theta_i^* = \arg\min_{\theta} ||\mathbf{f}(\theta, \mathbf{u}_i) - \mathbf{y}^0|| \Rightarrow \text{Optimisation.}$ 4: Compute  $\mathbf{J}_i$  with columns  $\partial \mathbf{f}(\theta_i^*, \mathbf{u}_i) / \partial \theta_k$ 5: Compute  $w_i = p(\theta_i^*) * (\det(\mathbf{J}_i^{\top}\mathbf{J}_i))^{-1/2}$ 6: Accept  $\theta_i^*$  as posterior sample with weight  $w_i$ .

As can readily be seen from Algorithm 1, ill-conditioned matrices  $\mathbf{J}_i^{\mathsf{T}} \mathbf{J}_i$  produce very large weights for the corresponding  $\boldsymbol{\theta}_i^*$ , possibly completely overshadowing the remaining samples and creating an approximate posterior density that is spiked at a single location (as in Figure 1). One may think that this issue can be easily fixed by regularising  $\mathbf{J}_i^{\mathsf{T}} \mathbf{J}_i$  before computing the determinant.

However, the issue goes deeper: ill-conditioned matrices  $\mathbf{J}_i^{\mathsf{T}} \mathbf{J}_i$  occur when a large parameter region around the optimum  $\boldsymbol{\theta}_i^*$  produces data with small distances. These regions are poorly approximated by point masses or infinitesimally small ellipsoids.<sup>1</sup> This means that if large parameter regions are a solution to  $||\mathbf{f}(\boldsymbol{\theta}, \mathbf{u}_i) - \mathbf{y}^0|| \le \epsilon$ , the likelihood function is nearly-constant around  $\boldsymbol{\theta}_i^*$  and cannot be approximated with a point mass, and amending the value of the weight cannot correct for this.

#### **3** Robust Optimisation Monte Carlo

We here develop a framework and concrete algorithms that have the benefits of OMC but do not collapse areas of similar likelihood into a point-mass.

#### 3.1 The Robust OMC Framework

We start from the basic characterisation of the finite sample version of the ABC posterior in Equation (3) which holds irrespective of OMC. Under this approximation, the expectation of an arbitrary function  $h(\theta)$  under the posterior  $p_{\epsilon}(\theta|\mathbf{x}^0)$  is

$$\mathbb{E}[h(\boldsymbol{\theta})|\mathbf{x}^{0}] = \frac{\int h(\boldsymbol{\theta})p(\boldsymbol{\theta})\frac{1}{n}\sum_{i=1}^{n}\mathbb{1}_{C_{\epsilon}}(\boldsymbol{\theta},\mathbf{u}_{i})\,d\boldsymbol{\theta}}{\int p(\boldsymbol{\theta})\frac{1}{n}\sum_{i=1}^{n}\mathbb{1}_{C_{\epsilon}}(\boldsymbol{\theta},\mathbf{u}_{i})\,d\boldsymbol{\theta}} \quad (4)$$
$$= \frac{\sum_{i=1}^{n}\int h(\boldsymbol{\theta})p(\boldsymbol{\theta})\mathbb{1}_{C_{\epsilon}^{i}}(\boldsymbol{\theta})\,d\boldsymbol{\theta}}{\sum_{i=1}^{n}\int p(\boldsymbol{\theta})\mathbb{1}_{C_{\epsilon}^{i}}(\boldsymbol{\theta})\,d\boldsymbol{\theta}}, \quad (5)$$

where  $C_{\epsilon}^{i} = \{\boldsymbol{\theta} : d(g(\boldsymbol{\theta}, \mathbf{u}_{i}), \mathbf{x}^{0}) \leq \epsilon\}$  is the set of parameters where, for a particular random seed or realisation of  $\mathbf{u}_{i}$ , the simulated data is within distance  $\epsilon$  from the observed data. Equation (5) features n integrals  $I_{i}$  in the numerator,

$$I_{i} = \int h(\boldsymbol{\theta}) p(\boldsymbol{\theta}) \mathbb{1}_{C_{\epsilon}^{i}}(\boldsymbol{\theta}) d\boldsymbol{\theta}, \qquad (6)$$

and similar ones in the denominator. The integrals are generally intractable but since they correspond to an expectation with respect to the prior  $p(\theta)$ , they could be approximated by a sample-based average because sampling from the prior is typically possible in likelihood-free inference problems. However, this would be inefficient in the case of a broad prior as most samples would give  $\mathbb{1}_{C_{\epsilon}^{i}}(\theta) = 0$ , i.e. they would essentially get rejected much like in rejection ABC. It is more efficient to sample from a proposal distribution  $q_{i}(\theta)$  that only has support on the acceptance region  $C_{\epsilon}^{i}$ . We will discuss how to construct such  $q_{i}(\theta)$  in Subsection 3.3. Assuming we have a

<sup>&</sup>lt;sup>1</sup>As pointed out in the original paper as a limitation of OMC, this happens if there are fewer summary statistics than parameters. But as shown here, this failure mode is more general. The proposed robust method provides a solution.

suitable  $q_i(\boldsymbol{\theta})$ , the integrals  $I_i$  can be approximated as

$$I_{i} = \int h(\boldsymbol{\theta}) \mathbb{1}_{C_{\epsilon}^{i}}(\boldsymbol{\theta}) \frac{p(\boldsymbol{\theta})}{q_{i}(\boldsymbol{\theta})} q_{i}(\boldsymbol{\theta}) d\boldsymbol{\theta}$$
(7)

$$\approx \frac{1}{m} \sum_{j=1}^{m} h(\boldsymbol{\theta}_{ij}) \mathbb{1}_{C_{\epsilon}^{i}}(\boldsymbol{\theta}_{ij}) \frac{p(\boldsymbol{\theta}_{ij})}{q_{i}(\boldsymbol{\theta}_{ij})}, \qquad (8)$$

where  $\theta_{ij} \sim q_i(\theta)$ , and equivalently for the integral in the denominator. Replacing the integrals in (5) with their sample-based approximations, we obtain

$$\mathbb{E}[h(\boldsymbol{\theta})|\mathbf{x}^{0}] \approx \frac{\sum_{i=1}^{n} \sum_{j=1}^{m} h(\boldsymbol{\theta}_{ij}) \mathbb{1}_{C_{\epsilon}^{i}}(\boldsymbol{\theta}_{ij}) \frac{p(\boldsymbol{\theta}_{ij})}{q_{i}(\boldsymbol{\theta}_{ij})}}{\sum_{i=1}^{n} \sum_{j=1}^{m} \mathbb{1}_{C_{\epsilon}^{i}}(\boldsymbol{\theta}_{ij}) \frac{p(\boldsymbol{\theta}_{ij})}{q_{i}(\boldsymbol{\theta}_{ij})}}.$$
 (9)

This expression corresponds to a weighted sample average of  $h(\theta)$ . Denoting  $\mathbb{E}[h(\theta)|\mathbf{x}^0]$  by  $\bar{h}$ , we have

$$\bar{h} \approx \frac{\sum_{ij} w_{ij} h(\boldsymbol{\theta}_{ij})}{\sum_{ij} w_{ij}}, \quad w_{ij} = \mathbb{1}_{C^i_{\epsilon}}(\boldsymbol{\theta}_{ij}) \frac{p(\boldsymbol{\theta}_{ij})}{q_i(\boldsymbol{\theta}_{ij})}, \quad (10)$$

where the  $w_{ij}$  are the (unnormalised) weights and the samples  $\theta_{ij}$  are drawn from  $q_i(\theta)$ . Since our test function  $h(\theta)$  has been arbitrary, this means that to obtain samples from the approximate posterior, we first draw n samples  $\mathbf{u}_i$ ,<sup>2</sup> thus defining the acceptance region  $C_{\epsilon}^i$ , and then m samples  $\theta_{ij}$  from the corresponding proposal distribution  $q_i(\theta)$ . This process is what we refer to as the Robust OMC approach.

Before discussing the construction of the proposal distributions  $q_i(\theta)$ , we show how OMC is obtained from (9) by making additional assumptions. Some of the assumptions can be easily violated in practice, which then leads to the failure mode pointed out above and illustrated in Figure 1.

#### 3.2 Connection to OMC

We present here the assumptions under which the proposed Robust OMC (ROMC) approach becomes standard OMC. It shows that ROMC is both more general and more robust than standard OMC.

**Theorem 3.1** Under the below assumptions, ROMC becomes equivalent to standard OMC.

Assumption 1. The distance  $d(g(\theta, \mathbf{u}), \mathbf{x}^0)$  is given by the Euclidean distance between summary statistics  $||\mathbf{f}(\theta, \mathbf{u}) - \mathbf{y}^0||$ .

Assumption 2. The proposal distribution  $q_i(\theta)$  is the uniform distribution on  $C_{\epsilon}^i$ .

Assumption 3. The acceptance regions  $C_{\epsilon}^{i}$  are approximated by the ellipsoid  $C_{\epsilon}^{i} = \{\boldsymbol{\theta} : (\boldsymbol{\theta} - \boldsymbol{\theta}_{i}^{*})^{\mathsf{T}} \mathbf{J}_{i}^{\mathsf{T}} \mathbf{J}_{i} (\boldsymbol{\theta} - \boldsymbol{\theta}_{i}^{*})^{\mathsf{T}} \mathbf{J}_{i}^{\mathsf{T}} \mathbf{J}_{i}^{\mathsf{T}} \mathbf{J}_{i} (\boldsymbol{\theta} - \boldsymbol{\theta}_{i}^{*})^{\mathsf{T}} \mathbf{J}_{i}^{\mathsf{T}} \mathbf{J}_{i} (\boldsymbol{\theta} - \boldsymbol{\theta}_{i}^{*})^{\mathsf{T}} \mathbf{J}_{i}^{\mathsf{T}} \mathbf{J}_{i}^{\mathsf{T}$ 

 $\theta_i^* \leq \epsilon$  where  $\mathbf{J}_i$  is the Jacobian matrix with columns  $\partial \mathbf{f}(\theta_i^*, \mathbf{u}_i) / \partial \theta_k$ .

Assumption 4. The matrix square root  $\mathbf{A}_i$  of  $\mathbf{J}_i^{\mathsf{T}} \mathbf{J}_i$  is full rank, i.e. rank $(\mathbf{A}_i) = dim(\boldsymbol{\theta})$ .

Assumption 5. The prior is constant on the acceptance regions  $C_{\epsilon}^{i}$ , i.e.  $p(\boldsymbol{\theta}) = p(\boldsymbol{\theta}_{i}^{*})$  for all *i*.

The proof of Theorem 3.1 is given in the appendix.

Assumptions 1 and 2 are of technical nature, but Assumption 1 highlights that ROMC can also use distances other than Euclidean ones. Assumption 3 and 4 show that OMC relies on  $C_{\epsilon}^i$  being well approximated by an ellipsoid of finite volume whose shape is determined by the local behaviour of  $f(\theta, u_i)$  at  $\theta_i^*$ . The failure case described in Subsection 2.2 and illustrated in Figure 1 is caused by a violation of these two assumptions. Assumption 5 is also important because it shows that e.g. strong smoothing of the empirical distribution defined by the weighted samples in OMC would ignore that the prior distribution may not be constant on the corresponding finite-sized ellipsoid.

#### 3.3 Robust OMC Algorithms

The Robust OMC (ROMC) framework has three key ingredients: the optimisation procedure as in OMC, the specification of the  $\epsilon$ -threshold as usual in ABC, and the proposal distribution. We here consider two sets of choices for these ingredients, resulting in Algorithms 2 and 3 detailed below. The former algorithm assumes access to (approximate) simulator gradients and can be run as post-processing to standard OMC, and the latter is gradient-free. The two algorithms show that the proposed ROMC framework is versatile and that it can be used to exploit specific properties of the model.

#### 3.3.1 Optimisation Step

To obtain the optimisation end point  $\theta_i^*$ , any optimisation algorithm can be used as long as it can minimise the distance with respect to  $\theta$ .

**Algorithm 2:** If gradients of the simulator are available, standard gradient-based optimisers are applicable.

Algorithm 3: If the simulator gradients are not available, we propose using Bayesian optimisation, which is a powerful optimisation scheme for objective functions that can be evaluated but whose gradients are not available (see e.g. Shahriari et al., 2016). In the simulations below, we use standard Bayesian optimisation (GPyOpt with the default settings) that builds a Gaussian Process surrogate model  $\hat{d}_i$  for each distance  $d(g(\theta, \mathbf{u}_i), \mathbf{x}^0)$  that needs to be minimised. The main purpose of the surro-

<sup>&</sup>lt;sup>2</sup>In practice, this is done by fixing the seeds of the simulator.

Algorithm 2 Boxed Robust OMC. Requires simulator gradients, possible as post-processing to standard OMC.

- 1: for  $i \leftarrow 1$  to n do
- 2: Obtain optimisation end point  $\theta_i^*$ .
- 3: Use curvature of  $\mathbf{J}_i \mathbf{J}_i$  to create a bounding box with volume V as described in Subsection 3.3.3.
- 4: Define a uniform distribution  $q_i(\theta)$  over the box.
- 5: for  $j \leftarrow 1$  to m do
- 6:  $\boldsymbol{\theta}_{ij} \sim q_i(\boldsymbol{\theta})$
- 7: Accept  $\boldsymbol{\theta}_{ij}$  as posterior sample with weight  $w_{ij} = \mathbb{1}_{C_{\epsilon}^{i}}(\boldsymbol{\theta}_{ij}) * p(\boldsymbol{\theta}_{ij}) * V$

gate model in Bayesian optimisation is to decide at which  $\theta$  to evaluate the distance next. This also applies to our situation but for ROMC, there are two further uses of the surrogate model: 1) it can be used to greatly speed up the acceptance check  $\mathbb{1}_{C_{\epsilon}^{i}}(\theta_{ij})$  in Equation (10) by using the surrogate distance rather than the true distance (see experiment 3); and 2) it can facilitate the construction of the proposal distribution  $q_i(\theta)$  as discussed below.

#### 3.3.2 Threshold

Algorithms 2 and 3: ROMC requires a value for the threshold  $\epsilon$  that occurs in the term  $\mathbb{1}_{C_{\epsilon}^{i}}(\boldsymbol{\theta}_{ij})$ . This requirement to set a threshold is similar to most ABC algorithms where it is typically chosen as a small quantile of the observed distances (for other solutions see e.g. the work by Beaumont et al. (2002); Blum and Francois (2010); Papamakarios and Murray (2016); Chen and Gutmann (2019); Simola et al. (2019)). We take a similar approach but base the value of  $\epsilon$  on the distances  $d_i^*$ at the optimisation end points,  $d_i^* = d(g(\boldsymbol{\theta}_i^*, \mathbf{u}_i), \mathbf{x}^0)$ . Since the  $d_i^*$  are the minimal distances obtained in the optimisation step, their values are much smaller than the distances that one would see in other ABC algorithms, and we can choose a large quantile. In our simulations we chose the 90% quantile of the  $d_i^*$  in order to be robust against bad optimisation instances. Since all  $\theta_{ij}$  are saved, the exact value for  $\epsilon$  can be changed later by the user without incurring any overheads.

#### 3.3.3 Proposal Distribution

We here describe two methods to construct the proposal distributions  $q_i(\theta)$ . We assume that the optimisation step has given us a sample  $\theta_i^*$  that is within  $C_{\epsilon}^i$ .

Algorithm 2: If simulator gradients are available, then it is possible to compute the matrix  $\mathbf{J}_i^{\mathsf{T}} \mathbf{J}_i$ . Its eigenvectors are orthogonal directions of highest curvature, along which we scan until we reach a point whose resulting distance no longer falls under the threshold. Doing so in each dimension specifies a box, and defining a uniform Algorithm 3 Ellipsoidal Robust OMC. Does not require simulator gradients.

- 1: for  $i \leftarrow 1$  to n do
- 2: Obtain optimisation end point  $\theta_i^*$  and GP model distance  $\hat{d}_i(\theta)$  using Bayesian optimisation.
- 3: Construct ellipse with volume V using  $\hat{d}_i(\theta)$  as described in Subsection 3.3.3.
- 4: Define a uniform distribution  $q_i(\theta)$  over ellipse.
- 5: **for**  $j \leftarrow 1$  to m **do**
- 6:  $\boldsymbol{\theta}_{ij} \sim q_i(\boldsymbol{\theta})$
- 7: Accept  $\boldsymbol{\theta}_{ij}$  as posterior sample with weight  $w_{ij} = \mathbb{1}_{C_{\epsilon}^{i}}(\boldsymbol{\theta}_{ij}) * p(\boldsymbol{\theta}_{ij}) * V$

distribution on this box gives us the proposal distribution  $q_i(\boldsymbol{\theta})$ . Since  $\mathbf{J}_i^{\mathsf{T}} \mathbf{J}_i$  is computed by standard OMC, this approach can be done entirely as post-processing to it.

Algorithm 3: We construct a box around the optimisation end point as above when we have simulator gradients, except that we use the Hessian of the GP model instead of  $\mathbf{J}_i^{\mathsf{T}} \mathbf{J}_i$ . To robustify the approach against e.g. inaccuracies in the GP model and hence estimation of the curvature, we sample parameter values from inside the box and compute their distance to the observed data using the posterior mean of the GP model. This incurs practically no overhead and is considerably faster than using the true distances if the simulator is expensive. These parameter-distance pairs are then used to train a quadratic regression model of the distance. The contour for which the distance is equal to a threshold defines an ellipsoid, and we use the uniform distribution on it as proposal distribution  $q_i(\boldsymbol{\theta})$ . Figure 2 visualises the construction and further details are provided in the appendix.



Figure 2: Algorithm 3, example proposal distribution  $q_i(\theta)$ . The contours show the GP model distance, the green dots visualise the true acceptance region  $C_{\epsilon}^i$ , and  $q_i(\theta)$  is the uniform distribution on the red ellipse which well approximates  $C_{\epsilon}^i$ .

#### 4 Experiments

We assess Robust OMC (ROMC) on three tasks and compare its performance to standard OMC. As reference, we use posteriors obtained by expensive Rejection ABC runs. In the appendix we further show comparison results to the exact posteriors when tractable. The accuracy is measured using the Jensen-Shannon divergence. We also contrast effective samples sizes, which are given by  $\text{ESS} = (\sum_{i=1}^{n} w_i)^2 / \sum_{i=1}^{n} w_i^2$ . Note that we did not perform additional comparisons to other ABC methods because such comparisons were already performed by Meeds and Welling (2015) with OMC showing clear improvements (about a factor of 10 fewer calls to the simulator per accepted sample). These advantages are inherited by ROMC.

We further compared the performance of ROMC to that of two simple heuristic fixes of OMC. In the first fix, we ignored a percentage of the smallest eigenvalues of the  $\mathbf{J}_i^{\mathsf{T}}\mathbf{J}_i$  matrices when computing the determinants, hence reducing the magnitude of the biggest weights. This is similar to using pseudo-determinants. For the second fix, we stabilised the  $\mathbf{J}_i^{\mathsf{T}}\mathbf{J}_i$  matrices by adding a constant value to the diagonals before computing the determinant. This stabiliser value was chosen by picking a given percentile of the magnitudes of the diagonal elements of all matrices.

#### 4.1 Experiment 1: ROMC Resolves OMC Failure

We first consider a simple simulator to illustrate that ROMC resolves the identified issue of standard OMC, namely that it collapses regions of similar likelihood into a single point. The simulator is defined such that the likelihood function is flat in the area around  $x^0$  and linear otherwise:

$$p(x|\theta) \sim \begin{cases} \theta^4 + u & \text{if } \theta \in [-0.5, 0.5] \\ \theta - c + u & \text{otherwise} \end{cases}$$
(11)

The parameter of interest is  $\theta$  and  $u \sim N(0, 1)$  is a nuisance parameter and the only source of randomness. The term  $c = (0.5 - 0.5^4)$  makes sure the function is continuous. Figure 8 in the appendix shows the simulator output for specific values of u. For posterior inference, we assume that the observed data is  $x^0 = 0$ . We used Algorithm 2 but exploited the fact that the box can be constructed analytically in this simple example.

Figure 1 in the introduction shows example posteriors. Despite generating samples which span the entire range of the prior, OMC assigns much higher weights to the samples in the middle of the flat region, resulting in a posterior that is overly confident at that point. Con-



Figure 3: Experiment 1. Comparing OMC, two heuristic fixes to OMC, and Robust OMC. Smaller divergences are better.

versely, ROMC can nearly perfectly reproduce the reference posterior. For OMC, the effective sample size divided by the number of total samples is  $\text{ESS}/n \approx 0.5$ , while for ROMC,  $\text{ESS}/n \approx 0.95$ .

Figure 3 shows how well OMC, ROMC, and the two discussed heuristic fixes of OMC can match the reference posterior at varying computational budgets. ROMC clearly outperforms OMC and both heuristics whatever the computation time. Additionally, we see that the two heuristic OMC methods perform similarly, with the pseudo-determinant version reaching a lower divergence. We will thus only consider that heuristic from now on.

#### 4.2 Experiment 2: ROMC as Post-processing

This experiment showcases that Robust OMC can be performed as post-processing to standard OMC. We assume that we can compute simulator gradients (which is required in OMC), and hence we use Algorithm 2. We consider the case were the summary statistics are not completely informative about the parameters, which is a scenario that comes up often when using ABC in real-world problems (e.g. Aeschbacher et al., 2012). As a prototypical example of this scenario, we infer the mean and variance of a normal distribution with only the sample average available as a summary statistic (see appendix for details). Since there is no direct information on the variance, the optimisation surfaces will be completely flat in one direction.

Figure 4 compares the methods for different run times against the reference Rejection ABC posterior. Heuristic OMC refers to the version based on pseudodeterminants, and we use the hyper-parameter value that produced the best result. As before, ROMC outperforms the alternatives. Figures of the posteriors themselves are shown in Figure 11 in the appendix, and Heuristic OMC performance as a function of its hyper-parameter is shown in Figure 13 in the appendix.



Figure 4: Experiment 2. Comparing OMC, heuristic OMC, and Robust OMC. Smaller divergences are better.



(a) Gray teapot with red light. (b) Red teapot with white light.

Figure 5: Example teapots (brightened for clarity). We use (a) as the observed data, but (b) is another possible explanation.

#### 4.3 Experiment 3: Gradient-free ROMC

Here, we do not make use of the simulator gradients, and use Algorithm 3. This example is about inverse-graphics. It involves a considerably more complicated simulator that takes as input a set of 20 parameters and deterministically renders an image of a object (in this case, a teapot) on a uniform background. This is based on the generative model used by Moreno et al. (2016). We focus on the task of learning the posterior distribution of two colour parameters in a setting where there are two possible explanations for the observed image and thus the posterior is expected to be bi-modal. The remaining 18 parameters are used as nuisance parameters. They control the illumination, shape, pose and other aspects of the objects (see appendix).

The five illumination parameters are the most relevant ones among the nuisance variables for the task considered. The first four parameters specify the global illumination strength, the directional light strength, and the directional light angle. The fifth one allows the directional lighting to be in one of two modes: either white or red. This is what causes the bi-modality of the posterior — if the observed image depicts a red teapot, it is both possible that it could be a grey teapot with red lighting, or a red teapot with white lighting (see Figure 5). We use the former case as the observed image in our experiments. The generative model was implemented using Open Differential Renderer (OpenDR, Loper and Black, 2014). While this renderer does allow us to compute the simulator gradients, we did not use them for any of our presented results. To compute the distance  $d(\mathbf{x}_i, \mathbf{x}^0)$  between a simulated image  $\mathbf{x}_i = g(\boldsymbol{\theta}, \mathbf{u}_i)$  and the observed image  $\mathbf{x}^0$ , we use a recognition model that predicts the target parameters  $\boldsymbol{\theta}$  for an input image  $\mathbf{x}$  and then compute the Euclidean distance between the two images' predicted parameters, so that

$$d(\mathbf{x}_i, \mathbf{x}^0) = \sqrt{\sum_{n=1}^{|\boldsymbol{\theta}|} (r_n(\mathbf{x}_i) - r_n(\mathbf{x}^0))^2}, \quad (12)$$

where  $r_n(\cdot)$  is the prediction of the *n*-th parameter. The parameter estimates can thus be viewed as summary statistics. We implemented the recognition model as a neural network that was pre-trained on data generated from the simulator using a broad prior on the nuisance variables under daylight (white illumination). For details about the neural network's architecture, training procedure, and performance, as well as for examples of the training data, see the appendix.

Obtaining a single optimisation end point  $\theta_i^*$ , i.e. minimising the distance  $d(\mathbf{x}_i, \mathbf{x}^0)$  with respect to the colour parameters, took approximately 2 minutes. We ran all simulations on a single computer only, and did not exploit the possibility to parallelise the inference. We based our posterior approximation on 250 optimisation end points. For ROMC, we generated 100 new samples per original optimisation end point.

The OMC and ROMC posteriors are shown in Figure 6, along with a Rejection ABC reference posterior and the predictions by the recognition network (see the appendix for results with the pseudo-determinant heuristic). First of all, we see that the recognition network prediction is off even though the network was well trained (see Figure 16 in the appendix). This is because red lightning conditions were not part of the training data and the recognition network does not well generalise towards this condition. Indeed, while still not accurate, the network favours the solution in Figure 5(b), which is closer to images typically seen during training. Among the Bayesian methods, OMC is overly confident at a single location, with an effective sample size of 1.2 (out of 250 total samples). What is more, its posterior is not centred on a viable solution. On the other hand, ROMC produces two distinct posterior modes that contain the two possible solutions as we would want in this scenario, and it generally matches the reference posterior. Remarkably, these results were obtained despite using a biased recognition network, which points to a general ability of ABC in dealing with systematic biases in recognition networks.



Figure 6: Experiment 3 posteriors. Left: Reference posterior, obtained after running Rejection ABC with a high number of samples. Middle: Standard OMC. ESS/ $n \approx 0.005$ . Right: ROMC with Algorithm 3: ESS/ $n \approx 0.97$ .



Figure 7: Experiment 3. Comparison between OMC and ROMC. For ROMC, we used Algorithm 3 without (red cross) and with (green square) GP acceleration.

Figure 7 compares the trade-off between accuracy and compute time for OMC (black) and ROMC (red and green). We see that ROMC as implemented in Algorithm 3 (red) is much more accurate than OMC but that it incurs an extra cost. This extra cost is due to the additional runs of the simulator that are needed for the acceptance check  $\mathbb{1}_{C_i^i}(\boldsymbol{\theta}_{ij})$  in step 7 of the algorithm. This cost could be reduced by parallelising the runs (which we did not do). Alternatively, we can use the GP model of the distance rather than the true distance in the acceptance check, which does not require additional runs of the simulator. We call this approach GP-ROMC. The figure shows that GP-ROMC (green) has almost the same accuracy as ROMC, but that it is much faster and only incurs a tiny overhead compared to OMC. The posterior for GP-ROMC is shown in Figure 12 in the appendix. In line with the numerical result, the posterior is very close to one obtained with exact acceptance checks.

#### **5** Conclusions

This paper dealt with the task of performing Bayesian inference for parametric models in the case where the likelihood is intractable but sampling from the model is possible. We considered Optimisation Monte Carlo (OMC) which has been shown to be a promising tool to efficiently sample from an approximate posterior. While efficient, we showed that OMC has the important shortcoming that it collapses regions of similar or nearconstant likelihood into a single point. This matters because OMC samples might thus severely under-represent the uncertainty in the posterior and hence produce overly confident predictions.

We addressed this issue by introducing the more general framework of Robust Optimisation Monte Carlo (ROMC) and two concrete algorithms implementing it. The ROMC framework can be considered to be a form of ABC where we use optimisation to automatically construct suitable and localised proposal distributions. The first algorithm can be run as a form of post-processing after standard OMC to correct for the identified pathology. The second algorithm, unlike OMC, can be used when (approximate) gradients are not available. It uses a surrogate model of the distance and we have seen that this approach can be used to almost entirely eliminate the extra cost of ROMC compared to OMC. It is hence reasonable to also use a surrogate model in the first algorithm if reducing compute cost is necessary.

We tested the proposed framework and algorithms on both prototypical toy examples and complex inference tasks from inverse-graphics, and found that the proposed ROMC approach did accurately estimate the posteriors while OMC did not.

#### Acknowledgements

This work was supported in part by the EPSRC Centre for Doctoral Training in Data Science, funded by the UK Engineering and Physical Sciences Research Council (grant EP/L016427/1) and the University of Edinburgh.

#### References

Simon Aeschbacher, Mark A Beaumont, and Andreas Futschik. A novel approach for choosing summary statistics in approximate Bayesian computation. *Genetics*, 192(3):1027–1047, 2012.

- Mark A Beaumont, Wenyang Zhang, and David J Balding. Approximate Bayesian computation in population genetics. *Genetics*, 162(4):2025–2035, 2002.
- E. Bernton, P.E. Jacob, M. Gerber, and C.P. Robert. Approximate Bayesian computation with the Wasserstein distance. *Journal of the Royal Statistical Society: Series B*, 2018.
- M. G. B. Blum, M. A. Nunes, D. Prangle, and S. A. Sisson. A comparative review of dimension reduction methods in approximate Bayesian computation. *Statistical Science*, 28(2):189–208, 2013.
- Michael Blum and Olivier Francois. Non-linear regression models for Approximate Bayesian Computation. *Statistics and Computing*, 20(1):63–73, 2010.
- Y. Chen and M.U. Gutmann. Adaptive Gaussian copula ABC. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AIS-TATS)*, 2019.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.
- Matteo Fasiolo, Simon N. Wood, Florian Hartig, and Mark V. Bravington. An extended empirical saddlepoint approximation for intractable likelihoods. *Electron. J. Statist.*, 12(1):1544–1578, 2018.
- Paul Fearnhead and Dennis Prangle. Constructing summary statistics for approximate Bayesian computation: semi-automatic approximate Bayesian computation. *Journal of the Royal Statistical Society: Series B* (*Statistical Methodology*), 74(3):419–474, 2012.
- Jean-Jacques Forneron and Serena Ng. A likelihood-free reverse sampler of the posterior distribution. In *Essays in Honor of Aman Ullah*, pages 389–415. Emerald Group Publishing Limited, 2016.
- Jean-Jacques Forneron and Serena Ng. The ABC of simulation estimation with auxiliary statistics. *Journal of Econometrics*, 205(1):112–139, 2018.
- Michael U Gutmann and Jukka Corander. Bayesian optimization for likelihood-free inference of simulatorbased statistical models. *The Journal of Machine Learning Research*, 17(1):4256–4302, 2016.
- Michael U Gutmann, Ritabrata Dutta, Samuel Kaski, and Jukka Corander. Likelihood-free inference via classification. *arXiv preprint arXiv:1407.4981*, 2014.
- Jarno Lintusaari, Michael U Gutmann, Ritabrata Dutta, Samuel Kaski, and Jukka Corander. Fundamentals and

recent developments in approximate Bayesian computation. *Systematic biology*, 66(1):e66–e82, 2017.

- Matthew M Loper and Michael J Black. Opendr: An approximate differentiable renderer. In *European Conference on Computer Vision*, pages 154–169. Springer, 2014.
- Pekka Marttinen, Nicholas J Croucher, Michael U Gutmann, Jukka Corander, and William P Hanage. Recombination produces coherent bacterial species clusters in both core and accessory genomes. *Microbial Genomics*, 1(5), 2015.
- Ted Meeds and Max Welling. Optimization Monte Carlo: Efficient and embarrassingly parallel likelihood-free inference. In *Advances in Neural Information Processing Systems*, pages 2080–2088, 2015.
- Pol Moreno, Christopher KI Williams, Charlie Nash, and Pushmeet Kohli. Overcoming occlusion with inverse graphics. In *Computer Vision–ECCV 2016 Work-shops*, pages 170–185. Springer, 2016.
- George Papamakarios and Iain Murray. Fast  $\varepsilon$ -free inference of simulation models with Bayesian conditional density estimation. In *Advances in Neural Information Processing Systems*, pages 1028–1036, 2016.
- L. F. Price, C. C. Drovandi, A. Lee, and D. J. Nott. Bayesian synthetic likelihood. *Journal of Computational and Graphical Statistics*, pages 1–11, March 2017.
- Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P. Adams, and Nando de Freitas. Taking the human out of the loop: A review of Bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2016.
- Umberto Simola, Jessica Cisewski-Kehe, Michael U. Gutmann, and Jukka Corander. Adaptive approximate Bayesian computation tolerance selection. *arXiv:1907.01505*, 2019.
- S.A. Sisson, Y Fan, and M.A. Beaumont. *Handbook of Approximate Bayesian Computation.*, chapter Overview of Approximate Bayesian Computation. Chapman and Hall/CRC Press, 2018.
- O. Thomas, R. Dutta, J. Corander, S. Kaski, and M.U. Gutmann. Likelihood-free inference by ratio estimation. arXiv:1611.10242, 2016.
- Tina Toni, David Welch, Natalja Strelkowa, Andreas Ipsen, and Michael PH Stumpf. Approximate Bayesian computation scheme for parameter inference and model selection in dynamical systems. *Journal of the Royal Society Interface*, 6(31):187–202, 2009.
- Simon N. Wood. Statistical inference for noisy nonlinear ecological dynamic systems. *Nature*, 466(7310): 1102–1104, 2010.

### Appendix

#### A Rejection ABC Algorithm

**Algorithm 4** Rejection ABC. Generates n independent samples  $\theta_i$  of the approximate posterior. Needs blackbox simulator  $g(\theta, \mathbf{u})$ , observed data  $\mathbf{x}^0$ , computational budget N, and number of accepted samples n.

- 1: for  $i \leftarrow 1$  to N do
- 2:  $\boldsymbol{\theta}_i \sim p(\boldsymbol{\theta})$  > Draw parameters  $\boldsymbol{\theta}_i$  from prior.
- 3:  $\mathbf{x}_i \sim g(\boldsymbol{\theta}_i, \cdot) \triangleright$  Simulate synthetic data using  $\boldsymbol{\theta}_i$ .
- 4:  $d_i = d(\mathbf{x}_i, \mathbf{x}^0) \triangleright$  Compute distance to  $\mathbf{x}^0$ .
- 5: Accept the *n* samples  $\theta_i$  with the lowest distance  $d_i$  as samples from the posterior.

A standard variant of Rejection ABC accepts only those samples whose distance falls under a threshold  $\epsilon$  specified by the user.

#### **B Proof of Theorem 3.1**

As a reminder, the assumptions under which ROMC becomes standard OMC are:

Assumption 1. The distance  $d(g(\theta, \mathbf{u}), \mathbf{x}^0)$  is given by the Euclidean distance between summary statistics  $||\mathbf{f}(\theta, \mathbf{u}) - \mathbf{y}^0||.$ 

Assumption 2. The proposal distribution  $q_i(\theta)$  is the uniform distribution on  $C_{\epsilon}^i$ .

Assumption 3. The acceptance regions  $C_{\epsilon}^{i}$  are approximated by the ellipsoid  $C_{\epsilon}^{i} = \{\boldsymbol{\theta} : (\boldsymbol{\theta} - \boldsymbol{\theta}_{i}^{*})^{\mathsf{T}} \mathbf{J}_{i}^{\mathsf{T}} \mathbf{J}_{i} (\boldsymbol{\theta} - \boldsymbol{\theta}_{i}^{*}) \leq \epsilon\}$  where  $\mathbf{J}_{i}$  is the Jacobian matrix with columns  $\partial \mathbf{f}(\boldsymbol{\theta}_{i}^{*}, \mathbf{u}_{i}) / \partial \theta_{k}$ .

Assumption 4. The matrix square root  $\mathbf{A}_i$  of  $\mathbf{J}_i^{\mathsf{T}} \mathbf{J}_i$  is full rank, i.e. rank $(\mathbf{A}_i) = \dim(\boldsymbol{\theta})$ .

Assumption 5. The prior is constant on the acceptance regions  $C_{\epsilon}^{i}$ , i.e.  $p(\boldsymbol{\theta}) = p(\boldsymbol{\theta}_{i}^{*})$  for all *i*.

Since settings of  $\mathbf{u}_i$  that result in empty sets  $C_{\epsilon}^i$  (i.e. no  $\theta_i^*$  exists such that the resulting distance is below  $\epsilon$ ) are excluded from affecting the approximate posterior (both in ROMC and in standard OMC), we here consider only the case of non-empty sets  $C_{\epsilon}^i$ .

We start from Equation (9) using for  $q_i(\boldsymbol{\theta})$  — as per Assumption 2 — the uniform distribution on  $C_{\epsilon}^i$  with density  $U_{\epsilon}^i(\boldsymbol{\theta})$ ,

$$U_{\epsilon}^{i}(\boldsymbol{\theta}) = \frac{1}{\operatorname{vol}(C_{\epsilon}^{i})} \mathbb{1}_{C_{\epsilon}^{i}}(\boldsymbol{\theta}).$$
(13)

Since the proposal distribution  $q_i(\boldsymbol{\theta})$  is zero outside  $C_{\epsilon}^i$ ,

we have  $\mathbb{1}_{C^i_{\epsilon}}(\boldsymbol{\theta}_{ij}) = 1$  for all  $\boldsymbol{\theta}_{ij}$  and hence

$$\mathbb{E}[h(\boldsymbol{\theta})|\mathbf{x}^{0}] \approx \frac{\sum_{i=1}^{n} \sum_{j=1}^{m} h(\boldsymbol{\theta}_{ij}) \operatorname{vol}(C_{\epsilon}^{i}) p(\boldsymbol{\theta}_{ij})}{\sum_{i=1}^{n} \sum_{j=1}^{m} \operatorname{vol}(C_{\epsilon}^{i}) p(\boldsymbol{\theta}_{ij})}.$$
 (14)

By Assumption 3,  $C_{\epsilon}^{i}$  is an ellipsoid with volume determined by the matrix square root **A** of  $\mathbf{J}_{i}^{\mathsf{T}}\mathbf{J}_{i}$ , as well as the value of  $\epsilon$ . It is possible to split the volume into a term determined by the shape of the ellipsoid and a term determined by  $\epsilon$ . With the change of variables  $\mathbf{w} = \mathbf{A}_{i}(\boldsymbol{\theta} - \boldsymbol{\theta}_{i}^{*})$ , we have:

$$\operatorname{vol}(C_{\epsilon}^{i}) = \int_{C_{\epsilon}^{i}} d\boldsymbol{\theta} = \int_{\mathbf{w}: ||\mathbf{w}||^{2} \le \epsilon} |\det(\mathbf{A}_{i})|^{-1} d\mathbf{w}$$
(15)
(15)

$$= |\det(\mathbf{A}_i)|^{-1} \operatorname{vol}(B_{\epsilon}), \tag{16}$$

where  $B_{\epsilon}$  denotes an  $\epsilon$ -ball in Euclidean space. By Assumption 4,  $|\det(\mathbf{A}_i)|^{-1}$  is finite. Low-rank matrices  $\mathbf{A}_i$  would correspond to ellipsoids that extend without bound into one (or more) directions. We thus obtain

$$\mathbb{E}[h(\boldsymbol{\theta})|\mathbf{x}^{0}] \approx \frac{\sum_{i=1}^{n} \sum_{j=1}^{m} h(\boldsymbol{\theta}_{ij}) \frac{\operatorname{vol}(B_{\epsilon})}{|\det(\mathbf{A}_{i})|} p(\boldsymbol{\theta}_{ij})}{\sum_{i=1}^{n} \sum_{j=1}^{m} \frac{\operatorname{vol}(B_{\epsilon})}{|\det(\mathbf{A}_{i})|} p(\boldsymbol{\theta}_{ij})} \approx \frac{\sum_{i=1}^{n} |\det(\mathbf{A}_{i})|^{-1} \sum_{j=1}^{m} h(\boldsymbol{\theta}_{ij}) p(\boldsymbol{\theta}_{ij})}{\sum_{i=1}^{n} |\det(\mathbf{A}_{i})|^{-1} \sum_{j=1}^{m} p(\boldsymbol{\theta}_{ij})}$$
(18)

where we cancelled  $\operatorname{vol}(B_{\epsilon})$  so that only the term  $|\det(\mathbf{A}_i)|^{-1}$  reflecting the geometry of the ellipsoid remains. Note that  $|\det(\mathbf{A}_i)|^{-1}$  can also be written as  $|\det(\mathbf{A}_i)|^{-1} = (\det \mathbf{J}_i^{\mathsf{T}}\mathbf{J}_i)^{-1/2}$ . By Assumption 5,  $p(\boldsymbol{\theta}_{ij}) = p(\boldsymbol{\theta}_i^*)$ , so that we have

$$\mathbb{E}[h(\boldsymbol{\theta})|\mathbf{x}^{0}] \approx \frac{\sum_{i=1}^{n} (\det \mathbf{J}_{i}^{\mathsf{T}} \mathbf{J}_{i})^{-\frac{1}{2}} p(\boldsymbol{\theta}_{i}^{*}) \sum_{j=1}^{m} h(\boldsymbol{\theta}_{ij})}{\sum_{i=1}^{n} (\det \mathbf{J}_{i}^{\mathsf{T}} \mathbf{J}_{i})^{-\frac{1}{2}} p(\boldsymbol{\theta}_{i}^{*})}$$
(19)

In this expression, the only dependency on  $\epsilon$  remains in the samples  $\theta_{ij} \sim U_{\epsilon}^{i}(\theta)$ . In the limit of infinitely small  $\epsilon$ ,  $U_{\epsilon}^{i}(\theta)$  becomes a Dirac delta distribution  $\delta(\theta - \theta_{i}^{*})$ centred at  $\theta_{i}^{*}$ . This means that the only possible sample from that distribution is  $\theta_{i}^{*}$  and hence that  $h(\theta_{ij}) = h(\theta_{i}^{*})$  for all j. In the limit of  $\epsilon \to 0$ , we thus obtain

$$\mathbb{E}[h(\boldsymbol{\theta})|\mathbf{x}^{0}] \approx \frac{\sum_{i=1}^{n} (\det \mathbf{J}_{i}^{\top} \mathbf{J}_{i})^{-\frac{1}{2}} p(\boldsymbol{\theta}_{i}^{*}) h(\boldsymbol{\theta}_{i}^{*})}{\sum_{i=1}^{n} (\det \mathbf{J}_{i}^{\top} \mathbf{J}_{i})^{-\frac{1}{2}} p(\boldsymbol{\theta}_{i}^{*})}.$$
 (20)

This expression is a weighted average using samples  $\theta_i^*$  and weights  $w_i$  as defined in Algorithm 1. This means that the stated assumptions yield the weighted posterior samples of OMC and concludes the proof.

#### C Constructing the Proposal Region for ROMC

Here we give more details on how we obtain the proposal distribution  $q_i(\boldsymbol{\theta})$ , using Figure 2 as a visual aid. We start with the optimisation end point  $\theta_i^*$  (blue cross in Figure 2). We also have the curvature matrix, which is  $\mathbf{J}_i^{\mathsf{T}} \mathbf{J}_i$  in Algorithm 2 or the Hessian of the GP model at  $\theta_i^*$  in Algorithm 3. Its eigenvectors are used to determine what we call scan directions - one per dimension, along with its opposite (orange lines). We move along these directions until the resulting distance to  $\mathbf{x}^0$  is no longer under a certain threshold. The distance is calculated with the GP model if it is available in order to speed this process up. These end points along the scan directions are then used to create a (loose) rectangular box, which is either the final proposal region (Algorithm 2), or is used to draw training data for the regression model to create an ellipse which is the final proposal region (Algorithm 3). In both cases, placing a uniform distribution on the proposal region gives the final proposal distribution  $q_i(\boldsymbol{\theta})$ .

While we do use  $\epsilon$  for the final  $\mathbb{1}_{C_{\epsilon}^{i}}(\boldsymbol{\theta}_{ij})$  check, it would be reasonable to make the proposal region slightly bigger in order to ensure we capture as much of the actual acceptance region as possible at the cost of rejecting a few more samples. We achieve this by specifying the threshold for finding  $q_i(\theta)$  to be bigger than the threshold used for  $\mathbb{1}_{C_{\epsilon}^{i}}(\boldsymbol{\theta}_{ij})$ . In Algorithm 2, we use  $\epsilon_{big}$  as defined by the 95% quantile of the optimisation end point distances to define the proposal region. In Algorithm 3, we use the same 95% quantile on the ellipse, and a bigger 97.5% quantile to define the loose box from which the training data for the regression model is drawn. It is important to note that just using the single 90% quantile for all these thresholds still produces a good final posterior, so the method is robust to that choice. The bigger thresholds we propose above produce a very slight performance improvement in practice (the divergence to the reference posterior with the bigger thresholds is about 1% smaller,<sup>3</sup> which indicates a better performance) and also make intuitive sense, which is why we use them in our final implementation.

For Algorithm 3, the regression model is quadratic and is trained via least squares. As it is quadratic, its contours are ellipsoidal. Thus, by finding the contour equal to  $\epsilon_{big}$ , we obtain an ellipse that is suitable for being the proposal region.

Additionally, to show that Algorithm 3 is reasonably robust to the construction of the loose box, we compare the above proposal construction method with an alternative one. This additional method works as follows: begin with a tiny box centred at the optimisation end point, aligned with the scan directions obtained from the GP model's Hessian. Sample  $\theta$  values uniformly from the box, compute their distance using the GP model, and check how many are under the acceptance threshold. Gradually expand this box until 50% of the sampled points are no longer within the threshold. That is the final loose region, which is then used to train the regression model and thus to produce an ellipse as before. The divergence between the resulting final posterior and the reference posterior is less than 3% bigger<sup>4</sup> than when we use Algorithm 3 as described previously. This difference is quite small, implying that different proposal region construction methods can work.

#### **D** Additional Information for Exp. 1

In this section, we further discuss Experiment 1 from Subsection 4.1. As a reminder, the likelihood function we use is defined by:

$$p(x|\theta) \sim \begin{cases} \theta^4 + u & \text{if } \theta \in [-0.5, 0.5] \\ \theta - c + u & \text{otherwise} \end{cases}$$
(21)

where  $u \sim N(0, 1)$  is a nuisance parameter and the only source of randomness, and the term  $c = (0.5 - 0.5^4)$  ensures the function is continuous. We assume that the observed data is  $x^0 = 0$  and that the distance is Euclidean.

Figure 8 shows the simulator output for specific values of u. In the case of u > 0, the simulator can never generate a data point that matches  $x^0 = 0$  for any  $\theta$ , although for u sufficiently close to 0, some  $\theta$  may result in a data point within the distance threshold.

For u < 0, we enter the interesting situation where there are two values for which the distance is 0, and the distance is non-zero between them. In other words, there are two possible zero-solutions to the OMC optimisation objective  $d(g(\theta, u_i), x^0)$ . This would imply that there can be two disjointed acceptance regions  $C_{\epsilon}^{i}$  for a single  $\theta_{i}^{*}$ if the threshold  $\epsilon$  is small enough. Currently, both OMC and our two Robust OMC implementations would not be able to capture the full disjointed region  $C^i_{\epsilon}$  in such a scenario - OMC would at best approximate only the size of the region around  $\theta_i^*$  with a weight, and Robust OMC would construct the box / ellipse and hence the proposal distribution only around  $\theta_i^*$  as well. This problem does not manifest in the results we have presented as we computed the correct acceptance region analytically. In general, this is a difficult issue to solve, although a simple fix

<sup>&</sup>lt;sup>3</sup>The exact numbers for the Jensen-Shannon divergences are 0.290 and 0.293.

<sup>&</sup>lt;sup>4</sup>The exact numbers for the Jensen-Shannon divergences are 0.290 and 0.298.



Figure 8: Experiment 1. Examples of the simulator's output and its distance to the observed data  $x^0 = 0$  for three specific values of the nuisance parameter u.



Figure 9: Experiment 1. Comparison of OMC, Robust OMC, and the two heuristic OMC methods against the true posterior. The reference Rejection ABC posterior's divergence is also shown.

(i.e. a potential improvement on Robust OMC) would be to restart the optimiser at different initial values in order to find all possible solutions to the optimisation objective. It is also possible that, with enough samples, errors from the disjointed regions would average out and thus the final posterior would still be correct.

In the main text, Figure 3 compared OMC, Robust OMC, and two heuristic OMC methods against a reference posterior obtained from an expensive Rejection ABC run. Here we show the comparison made against the true posterior (which can be computed analytically) in Figure 9. As before, Robust OMC outperforms the other methods. It does not quite reach the level of the reference Rejection ABC method (blue dashed line) but this is to be expected as the reference was ran for a much longer time than the other methods.

#### E Additional Information for Exp. 2

Here we present further details and results for Experiment 2 discussed in Subsection 4.2. Recall that the task was about inferring the parameters of a 2D Gaus-



Figure 10: Experiment 2. Comparison of OMC and Robust OMC against the true posterior. The reference posterior's divergence is also shown.

sian with a non-informative summary statistic, namely just the mean of a sample from the Gaussian, while also having access to the simulator gradients.

We assume that the sample size is M = 25, and that the observed sample average is  $\mu^0 = 1$ . We chose a Gaussian prior for the mean and Inverse-Gamma prior for the standard deviation:  $p(\mu) = \mathcal{N}(0,5), p(\sigma) =$ Inv-Gamma(0.2, 1). To compute the exact posterior, we used the fact that for a sample of size M from a normal  $\mathcal{N}(\mu, \sigma^2)$  the sample mean is distributed according to  $\mathcal{N}(\mu, \sigma^2/M)$ .

In Figure 10, we show the performance of the methods compared against the true posterior rather than the reference ABC one used in the main text. The Jensen-Shannon divergence between Robust OMC's posterior and the true posterior is much lower than for the other methods. The actual posteriors themselves are shown in Figure 11. OMC correctly identifies the marginal over  $\mu$  but fails to do so for  $\sigma$  and does not match the reference posterior. On the other hand, Robust OMC reasonably matches the reference. Also note that OMC's effective sample size divided by the number of total samples is ESS/ $n \approx 0.01$ , implying that the vast majority of the samples are ignored. Conversely, the corresponding value for Robust OMC is ESS/ $n \approx 0.55$ , which is a significant improvement.

Figure 13 additionally compares the performance of our pseudo-determinant heuristic fix to OMC for different values of the hyper-parameter. This hyper-parameter represents the number of eigenvalues ignored when computing the determinants of the  $\mathbf{J}_i^{\mathsf{T}} \mathbf{J}_i$  matrices and hence the weights. While the divergence to the reference posterior does change, there is still a large gap between the best heuristic OMC and the robust OMC result.

#### F Additional Information for Exp. 3

Additional results. Figure 12 qualitatively compares GP-ROMC — the approach where we speed up the final distance check in Algorithm 3 by using the GP model instead of the simulator - to the standard Robust OMC approach and the reference posterior. Additionally, Figure 14 compares overall Robust OMC performance to that of OMC and Heuristic OMC. Comparing Robust OMC to Heuristic OMC, we noticed that as more eigenvalues are ignored for Heuristic OMC, the weights become more similar and the accuracy of the posterior improves. This is because, in this particular example and unlike before, the unweighted samples  $\theta_i^*$  do reasonably represent the posterior so that setting all weights to a constant provides a reasonable solution. However, such tuning is not possible in practice where a reference posterior is not available.

All renderer parameters. The full list of parameters we use for the renderer in Experiment 3 in Subsection 4.3 is as follows:

- Ten shape parameters. The object's exact shape is based on a morphable mesh specified by Principal Component Analysis. The 10 dimensions used are the 10 highest principal components.
- Two rotation parameters, specifically the azimuth and elevation. The camera is always centred at the midpoint of the object.
- Three colour parameters an RGB array which globally identifies the colour of the object. The first two (the red and green channels) are the target parameters  $\theta$  over which we perform inference in Experiment 3.
- Five illumination parameters that characterise the lighting on the object. Unlike Moreno et al. (2016) who use spherical harmonics to model illumina-

tion, we use single-source directional lighting as it is more intuitive and natural.

Recognition network details. The network we used has 3 convolutional layers, each with 64 5x5 filters and 2x2 max pooling, followed by 2 linear layers with 256 and 64 hidden units respectively. Each layer uses ReLU activation functions except the final layer which uses an identity activation. The network parameters were learned with Adagrad (Duchi et al., 2011) as it showed the most robustness to the values of the hyper-parameters of the neural network training procedure (batch size, learning rate, and dropout probability), which in turn were chosen via hyper-parameter optimisation. Additionally, Figure 15 shows samples from the training set used to train the recognition network. Note that there is a reasonable amount of variability in shape, pose, illumination, and colour. Figure 16 shows that the learned recognition network is reasonably good at reconstructing the parameters for test images from the training data.



Figure 11: Experiment 2 posteriors. For OMC,  $\text{ESS}/n \approx 0.02$ , implying that the vast majority of the samples are ignored. For Robust OMC,  $\text{ESS}/n \approx 0.55$ , which is a significant improvement.



Figure 12: Experiment 3 posteriors. Similar to Figure 6, but with the GP-ROMC approach (middle) shown for comparison as well.



Figure 13: Performance of the heuristic pseudodeterminant fix to OMC as a function of its hyperparameter in Experiment 2, compared against Robust OMC ran with roughly the same computational budget.



Figure 14: Performance of the heuristic pseudodeterminant fix to OMC in Experiment 3, compared against OMC and Robust OMC.



Figure 15: Experiment 3. Examples from the training set used to train the recognition network.



Figure 16: Experiment 3. Neural network predictions for the first two colour parameters  $c_0$ , and  $c_1$ . Line where true value is exactly equal to the predicted value is given for reference.