

Adapting the ABC Distance Function

Dennis Prangle^{*†}

Abstract. Approximate Bayesian computation performs approximate inference for models where likelihood computations are expensive or impossible. Instead simulations from the model are performed for various parameter values and accepted if they are close enough to the observations. There has been much progress on deciding which summary statistics of the data should be used to judge closeness, but less work on how to weight them. Typically weights are chosen at the start of the algorithm which normalise the summary statistics to vary on similar scales. However these may not be appropriate in iterative ABC algorithms, where the distribution from which the parameters are proposed is updated. This can substantially alter the resulting distribution of summary statistics, so that different weights are needed for normalisation. This paper presents two iterative ABC algorithms which adaptively update their weights and demonstrates improved results on test applications.

Keywords: likelihood-free inference, population Monte Carlo, quantile distributions, Lotka–Volterra.

1 Introduction

Approximate Bayesian computation (ABC) is a family of approximate inference methods which can be used when the likelihood function is expensive or impossible to compute but simulation from the model is straightforward. The simplest algorithm is a form of rejection sampling. Here parameter values are drawn from the prior distribution and corresponding datasets simulated. Each simulation is converted to a vector of summary statistics $\mathbf{s} = (s_1, s_2, \dots, s_m)$ and a distance between this and the summary statistics of the observed data, \mathbf{s}_{obs} , is calculated. Parameters producing distances below some threshold are accepted and form a sample from an approximation to the posterior distribution.

The choice of summary statistics has long been recognised as being crucial to the quality of the approximation (Beaumont et al., 2002), but there has been less work on the role of the distance function. A popular distance function is weighted Euclidean distance:

$$d(\mathbf{s}, \mathbf{s}_{\text{obs}}) = \left[\sum_{i=1}^m \left(\frac{s_i - s_{\text{obs},i}}{\sigma_i} \right)^2 \right]^{1/2} \quad (1)$$

where σ_i is an estimate of the prior predictive standard deviation of the i th summary

^{*}Department of Mathematics and Statistics, Newcastle University, UK, dennis.prangle@gmail.com

[†]Thanks to Michael Stumpf and Scott Sisson for helpful discussions, and three anonymous referees for feedback including suggesting Algorithm 4. The main part of this work was completed while the author was supported by a Richard Rado postdoctoral fellowship from the University of Reading.

statistic. In ABC rejection sampling a convenient estimate is the empirical standard deviation of the simulated s_i values. Scaling by σ_i in (1) normalises the summaries so that they vary over roughly the same scale, preventing the distance being dominated by the most variable summary.

This paper concerns the choice of distance in more efficient iterative ABC algorithms, in particular those of Toni et al. (2009), Sisson et al. (2009) and Beaumont et al. (2009). The first iteration of these algorithms is the ABC rejection sampling algorithm outlined above. The sample of accepted parameters is used to construct an importance density. An ABC version of importance sampling is then performed. This is similar to ABC rejection sampling, except parameters are sampled from the importance density rather than the prior, and the output sample is weighted appropriately to take this change into account. The idea is to concentrate computational resources on performing simulations for parameter values likely to produce good matches. The output of this step is used to produce a new importance density and perform another iteration, and so on. In each iteration the acceptance threshold is reduced, resulting in increasingly accurate approximations. Full details of the Toni et al. (2009) implementation are reviewed later.

Weighted Euclidean distance is commonly used in these algorithms with σ_i values determined in the first iteration. However there is no guarantee that these will normalise the summary statistics produced in later iterations, as these are no longer drawn from the prior predictive. This paper proposes two variant iterative ABC algorithms which update their σ_i values to appropriate values at each iteration. It is demonstrated that these algorithms provide substantial advantages in applications. Also, they do not require any extra simulations to be performed solely for tuning. Therefore even when a non-adaptive distance performs adequately, there is no major penalty in using the new approach. (Some additional calculations are required – calculating more σ_i values and more expensive distance calculations – but these form a negligible part of the overall computational cost.)

One of the proposed algorithms has similarities to the iterative ABC methods of Sedki et al. (2012) and Bonassi and West (2015). These postpone deciding some elements of the tuning of iteration t until during that iteration. Algorithm 5 also uses this strategy but for different tuning decisions: the distance function and the acceptance threshold. Another related paper is Fasiolo and Wood (2015) which contains an illustration of the difficulty of choosing ABC distance weights non-adaptively.

The remainder of the paper is structured as follows. Section 2 reviews ABC algorithms. This includes some novel material on the convergence of iterative ABC methods. Full technical details of these convergence results are given in supplementary material (Prangle, 2016). Section 3 discusses weighting summary statistics in a particular ABC distance function. Section 4 details the proposed algorithms. Several examples are given in Section 5. Finally, Section 6 summarises the work and discusses potential extensions. Computer code to implement the methods of this paper in the Julia programming language (Bezanson et al., 2012) is available at <https://github.com/dennisprangle/ABCDistances.jl>.

2 Approximate Bayesian computation

This section sets out the necessary background on ABC algorithms. Several review papers (e.g. Beaumont, 2010; Csilléry et al., 2010; Marin et al., 2012) give detailed descriptions of other aspects of ABC, including tuning choices and further algorithms. Sections 2.1 and 2.2 review ABC versions of rejection sampling and population Monte Carlo (PMC). Section 2.3 contains novel material on the convergence of ABC algorithms.

2.1 ABC rejection sampling

Consider Bayesian inference for parameter vector θ under a model with density $\pi(\mathbf{y}|\theta)$. Let $\pi(\theta)$ be the prior density and \mathbf{y}_{obs} represent the observed data. It is assumed that $\pi(\mathbf{y}|\theta)$ cannot easily be evaluated but that it is straightforward to sample from the model. ABC rejection sampling (Algorithm 1) exploits this to sample from an approximation to the posterior density $\pi(\theta|\mathbf{y})$. It requires several tuning choices: number of simulations N , a threshold $h \geq 0$, a function $S(\mathbf{y})$ mapping data to a vector of summary statistics, and a distance function $d(\cdot, \cdot)$.

Algorithm 1 ABC-rejection

1. Sample θ_i^* from $\pi(\theta)$ independently for $1 \leq i \leq N$.
 2. Sample \mathbf{y}_i^* from $\pi(\mathbf{y}|\theta_i^*)$ independently for $1 \leq i \leq N$.
 3. Calculate $\mathbf{s}_i^* = S(\mathbf{y}_i^*)$ for $1 \leq i \leq N$.
 4. Calculate $d_i^* = d(\mathbf{s}_i^*, \mathbf{s}_{\text{obs}})$ (where $\mathbf{s}_{\text{obs}} = S(\mathbf{y}_{\text{obs}})$).
 5. Return $\{\theta_i^* | d_i^* \leq h\}$.
-

The threshold h may be specified in advance. Alternatively it can be calculated following step 4. For example a common choice is to specify an integer k and take h to be the k th smallest of the d_i^* values (Biau et al., 2015).

2.2 ABC-PMC

Algorithm 2 is an iterative ABC algorithm taken from Toni et al. (2009). Very similar algorithms were also proposed by Sisson et al. (2009) and Beaumont et al. (2009). The latter note that this approach is an ABC version of population Monte Carlo (Cappé et al., 2004), so it is referred to here as ABC-PMC. The algorithm involves a sequence of thresholds, $(h_t)_{t \geq 1}$. Similarly to h in ABC-rejection, this can be specified in advance or during the algorithm, as discussed below.

The algorithm samples parameters from the importance density

$$q_t(\theta) = \begin{cases} \pi(\theta) & \text{if } t = 1, \text{ or } t = 2 \text{ and } h_1 = \infty \quad (2a) \\ \sum_{i=1}^N w_i^{t-1} K_t(\theta|\theta_i^{t-1}) / \sum_{i=1}^N w_i^{t-1} & \text{otherwise.} \quad (2b) \end{cases}$$

In the first iteration (and sometimes the second, as discussed shortly) $q_t(\theta)$ is the prior. Otherwise (2b) is used, which effectively samples from the previous weighted population and perturbs the result using kernel K_t . Beaumont et al. (2009) show that a good choice

Algorithm 2 ABC-PMC (with the option of adaptive h_t)

Initialisation

1. Let $t = 1$.

Main loop

2. Repeat following steps until there are N acceptances.
 - (a) Sample θ^* from importance density $q_t(\theta)$ given in equation (2).
 - (b) If $\pi(\theta^*) = 0$ reject and return to (a).
 - (c) Sample \mathbf{y}^* from $\pi(\mathbf{y}|\theta_i^*)$ and calculate $\mathbf{s}^* = S(\mathbf{y}^*)$.
 - (d) Accept if $d(\mathbf{s}^*, \mathbf{s}_{\text{obs}}) \leq h_t$.

Denote the accepted parameters as $\theta_1^t, \dots, \theta_N^t$ and the corresponding distances as d_1^t, \dots, d_N^t .

3. Let $w_i^t = \pi(\theta_i^t)/q_t(\theta_i^t)$ for $1 \leq i \leq N$.
4. (Optional) Let h_{t+1} be the α quantile of the d_i^t values.
5. Increment t to $t + 1$.

End of loop

of the latter is

$$K_t(\theta|\theta') = \phi(\theta', 2\Sigma_{t-1}),$$

where ϕ is the density of a normal distribution and Σ_{t-1} is the empirical variance matrix of $(\theta_i^{t-1})_{1 \leq i \leq N}$ calculated using weights $(w_i^{t-1})_{1 \leq i \leq N}$

As mentioned above, the sequence of thresholds can be specified in advance. However it is hard to do this well. A popular alternative (Drovandi and Pettitt, 2011a) is to choose the thresholds adaptively by setting h_t at the end of iteration $t - 1$ to be the α quantile of the accepted distances (n.b. $\alpha < 1$ is assumed throughout the paper). An optional step, step 4, is included in Algorithm 2 to implement this method. Alternative updating rules for h_t have been proposed such as choosing it to reduce an estimate of effective sample size by a prespecified proportion (Del Moral et al., 2012) or using properties of the predicted ABC acceptance rate (Silk et al., 2013).

If step 4 is used this leaves h_1 and α as tuning choices. A simple default for h_1 is ∞ , in which case all simulations are accepted when $t = 1$. In this case (2b) would give $q_2(\theta)$ as simply a modified prior with inflated variance, which is not a sensible importance density. Therefore (2) takes $q_2(\theta) = \pi(\theta)$ in this case. This is a minor novelty of this presentation of the algorithm.

A practical implementation of Algorithm 2 requires a condition for when to terminate. In this paper the total number of datasets to simulate is specified as a tuning parameter and the algorithm stops once a further simulation is required. Some alternative are possible, such as stopping once the algorithm falls below a target value for h_t or the acceptance rate.

Several variations on Algorithm 2 have been proposed which are briefly discussed in Section 6. Some of these are ABC versions of sequential Monte Carlo (SMC). The phrase “iterative ABC” will be used to cover ABC-PMC and ABC-SMC.

2.3 Convergence of ABC-PMC

Conditions C1–C5 ensure that Algorithm 2 converges on the posterior density in an appropriate sense as the number of iterations tends to infinity. This follows from Theorem 1 which is detailed in the supplementary material. Although only finite computational budgets are available in practice, such convergence at least guarantees that the target distribution become arbitrarily accurate as computational resources are increased.

- C1. $\theta \in \mathbb{R}^n$, $\mathbf{s} \in \mathbb{R}^m$ for some m, n and these random variables have density $\pi(\theta, \mathbf{s})$ with respect to Lebesgue measure.
- C2. The sets $A_t = \{\mathbf{s} | d(\mathbf{s}, \mathbf{s}_{\text{obs}}) \leq h_t\}$ are Lebesgue measurable.
- C3. $\pi(\mathbf{s}_{\text{obs}}) > 0$.
- C4. $\lim_{t \rightarrow \infty} |A_t| = 0$ (where $|\cdot|$ represents Lebesgue measure).
- C5. The sets A_t have *bounded eccentricity*.

Bounded eccentricity is defined in the supplementary material. Roughly speaking, it requires that under any projection of A_t to a lower dimensional space the measure still converges to zero.

Condition C1 is quite strong, ruling out discrete parameters and summary statistics, but makes proof of Theorem 1 straightforward. Condition C2 is a mild technical requirement. The other conditions provide insight into conditions required for convergence. Condition C3 requires that it must be possible to simulate \mathbf{s}_{obs} under the model. Condition C4 requires that the acceptance regions A_t shrink to zero measure. For most distance functions this corresponds to $\lim_{t \rightarrow \infty} h_t = 0$. It is possible for this to fail. Some examples encountered by the author in practice follow. One is when datasets close to \mathbf{s}_{obs} cannot be produced under the model of interest. Alternatively, even if \mathbf{s}_{obs} can occur under the model, the algorithm may converge on importance densities on θ under which it is impossible. This corresponds to concentrating on the wrong mode of the ABC target distribution in an early iteration. Finally, condition C5 prevents A_t converging to a set where some but not all summary statistics are perfectly matched.

Conditions C4 and C5 can be used to check which distance functions are sensible to use in ABC-PMC, usually by investigating whether they hold when $h_t \rightarrow 0$. For example it is straightforward to show this is the case when $d(\cdot, \cdot)$ is a metric induced by a norm.

3 Weighted Euclidean distance in ABC

This paper concentrates on using weighted Euclidean distance in ABC. Section 3.1 discusses this distance and how to choose its weights. Section 3.2 illustrates its usefulness in a simple example.

3.1 Definition and usage

Consider the following distance:

$$d(\mathbf{x}, \mathbf{y}) = \left[\sum_{i=1}^m \{\omega_i(x_i - y_i)\}^2 \right]^{1/2}. \quad (3)$$

If $\omega_i = 1$ for all i , this is *Euclidean distance*. Otherwise it is a form of *weighted Euclidean distance*.

Many other distance functions can be used in ABC, as discussed in Section 2.3, for example weighted L_1 distance $d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^m \omega_i |x_i - y_i|$. To the author's knowledge the only published comparison of distance functions is by McKinley et al. (2009), which found little difference between the alternatives. Owen et al. (2015) report the same conclusion but not the details. This finding is also supported in unpublished work by the author of this paper. Given these empirical results this paper focuses on (3) as it is a simple choice, but no claims are made for its optimality. Some further discussion on this is given in Section 6.

Summary statistics used in ABC may vary on substantially different scales. In the extreme case Euclidean distance will be dominated by the most variable. To avoid this, weighted Euclidean distance is generally used. This usually takes $\omega_i = 1/\sigma_i$ where σ_i is an estimate of the scale of the i th summary statistic. (Using this choice in weighted Euclidean distance gives the distance function (1) discussed in the introduction.)

A popular choice (e.g. Beaumont et al., 2002) of σ_i is the empirical standard deviation of the i th summary statistic under the prior predictive distribution. Csilléry et al. (2012) suggest using median absolute deviation (MAD) instead since it is more robust to large outliers. MAD is used throughout this paper. For many ABC algorithms these σ_i values can be calculated without requiring any extra simulations. For example this can be done between steps 3 and 4 of ABC-rejection. ABC-PMC can be modified similarly, resulting in Algorithm 3, which also updates h_t adaptively. (n.b. All of the ABC-PMC convergence discussion in Section 2.3 also applies to this modification.)

3.2 Illustration

As an illustration, Figure 1 shows the difference between using Euclidean and weighted Euclidean distance with $\omega_i = 1/\sigma_i$ within ABC-rejection. Here σ_i is calculated using MAD. For both distances the acceptance threshold is tuned to accept half the simulations. In this example Euclidean distance mainly rejects simulations where s_1 is far from its observed value: it is dominated by this summary. Weighted Euclidean distance also rejects simulations where s_2 is far from its observed value and is less stringent about s_1 .

Which of these distances is preferable depends on the relationship between the summaries and the parameters. For example if s_1 were the only informative summary, then Euclidean distance would be preferable. In practice, this relationship may not be known. Weighted Euclidean distance is then a sensible choice as both summary statistics contribute to the acceptance decision.

Algorithm 3 ABC-PMC with adaptive h_t and $d(\cdot, \cdot)$ **Initialisation**

1. Let $t = 1$ and $h_1 = \infty$.

Main loop

2. Repeat following steps until there are N acceptances.
 - (a) Sample θ^* from importance density $q_t(\theta)$ given in equation (2).
 - (b) If $\pi(\theta^*) = 0$ reject and return to (a).
 - (c) Sample \mathbf{y}^* from $\pi(\mathbf{y}|\theta_i^*)$ and calculate $\mathbf{s}^* = S(\mathbf{y}^*)$.
 - (d) Accept if $d(\mathbf{s}^*, \mathbf{s}_{\text{obs}}) \leq h_t$ (if $t = 1$ always accept).
3. If $t = 1$:
 - (a) Calculate $(\sigma_1, \sigma_2, \dots)$, a vector of MADs for each summary statistic, calculated from all the simulations in step 2 (including those rejected).
 - (b) Define $d(\cdot, \cdot)$ as the distance (3) using weights $(\omega_i)_{1 \leq i \leq m}$ where $\omega_i = 1/\sigma_i$.
 Denote the accepted parameters as $\theta_1^t, \dots, \theta_N^t$ and the corresponding distances as d_1^t, \dots, d_N^t .
4. Let $w_i^t = \pi(\theta_i^t)/q_t(\theta_i^t)$ for $1 \leq i \leq N$.
5. Let h_t be the α quantile of the d_i^t values.
6. Increment t to $t + 1$.

End of loop

This heuristic argument supports the use of weighted Euclidean distance in ABC more generally. One particular case is when low dimensional informative summary statistics have been selected, for example by the methods reviewed in Blum et al. (2013). In this situation all summaries are known to be informative and should contribute to the acceptance decision.

Note that in Figure 1 the observed summaries \mathbf{s}_{obs} lie close to the centre of the set of simulations. When some observed summaries are hard to match by model simulations this is not the case. ABC distances could now be dominated by the summaries which are hardest to match. How to weight summaries in this situation is discussed in Section 6.

4 Methods: Iterative ABC with an adaptive distance

The previous section discussed normalising ABC summary statistics using estimates of their scale under the prior predictive distribution. This prevents any summary statistic dominating the acceptance decision in ABC-rejection or the first iteration of Algorithm 3, where the simulations are generated from the prior predictive. However in later iterations of Algorithm 3 the simulations may be generated from a very different distribution so that this scaling is no longer appropriate. This section presents two versions of ABC-PMC which avoid this problem by updating the distance function at each iteration. Normalisation is now based on the distribution of summary statistics generated in

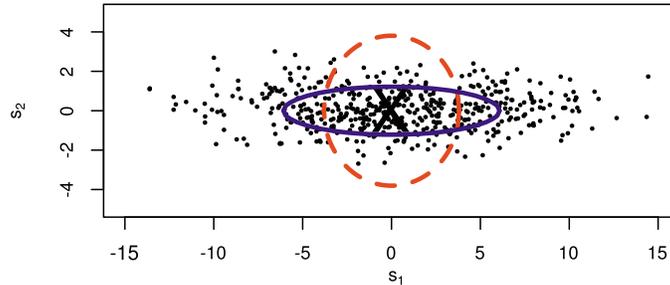


Figure 1: An illustration of distance functions in ABC rejection sampling. The points show simulated summary statistics s_1 and s_2 . The observed summary statistics are taken to be $(0, 0)$ (black cross). Acceptance regions are shown for two distance functions, Euclidean (red dashed circle) and weighted Euclidean with MAD reciprocals as weights (blue solid ellipse). These show the sets within which summaries are accepted. The acceptance thresholds have been tuned so that each region contains half the points.

the previous (Algorithm 4) or current (Algorithm 5) iteration. The proposed algorithms are presented in Sections 4.1 and 4.2.

An approach along these lines has the danger that the summary statistic acceptance regions at each iteration no longer form a nested sequence of subsets converging on the point $\mathbf{s} = \mathbf{s}_{\text{obs}}$. To avoid this, the proposed algorithms only accept a simulated dataset at iteration t if it also meets the acceptance criteria of *every previous iteration*. This can be viewed as sometimes modifying the t th distance function to take into account information from previous iterations. Section 4.3 discusses convergence in more depth.

4.1 First proposed algorithm

Algorithm 4 is a straightforward modification of Algorithm 3 which updates its distance function at each iteration using scales derived from the previous iteration's simulations. The first iteration accepts everything so no distance function is required. This acts as an initial tuning step. Note that scales are based on both accepted and rejected simulations from the previous iteration. This is because using just the accepted simulations would mean the scales are sometimes mainly determined by the previous acceptance rule, restricting the scope for adaptation.

Storing all simulated \mathbf{s}^* vectors to calculate scale estimates in step 3 of Algorithm 4 can be impractical. In practice storage is stopped after the first few thousand simulations, and scale estimation is done using this subset. Other tuning details of Algorithm 4 – the choice of perturbation kernel K_t and the rule to terminate the algorithm – are implemented as described earlier for ABC-PMC.

4.2 Second proposed algorithm

Algorithm 4 normalises simulations in iteration t based on scales derived in the preceding iteration. This could be inappropriate if two consecutive iterations sometimes

Algorithm 4 ABC-PMC with adaptive h_t and $d^t(\cdot, \cdot)$ **Initialisation**

1. Let $t = 1$ and $h_1 = \infty$.

Main loop

2. Repeat following steps until there are N acceptances.
 - (a) Sample θ^* from importance density $q_t(\theta)$ given in equation (2).
 - (b) If $\pi(\theta^*) = 0$ reject and return to (a).
 - (c) Sample \mathbf{y}^* from $\pi(\mathbf{y}|\theta_i^*)$ and calculate $\mathbf{s}^* = S(\mathbf{y}^*)$.
 - (d) If $t = 1$ accept. Otherwise accept if $d^i(\mathbf{s}^*, \mathbf{s}_{\text{obs}}) \leq h_i$ for all $2 \leq i \leq t$.
3. Calculate $(\sigma_1^t, \sigma_2^t, \dots)$, a vector of MADs for each summary statistic, calculated from all the simulations in step 2 (including those rejected).
4. Define $d^{t+1}(\cdot, \cdot)$ as the distance (3) using weights $(\omega_i)_{1 \leq i \leq m}$ where $\omega_i = 1/\sigma_i$. Denote the accepted parameters as $\theta_1^t, \dots, \theta_N^t$ and the corresponding distances under $d^{t+1}(\cdot, \cdot)$ as $d_1^{t+1}, \dots, d_N^{t+1}$.
5. Let $w_i^t = \pi(\theta_i^t)/q_t(\theta_i^t)$ for $1 \leq i \leq N$.
6. Let h_{t+1} be the α quantile of the d_i^{t+1} values.
7. Increment t to $t + 1$.

End of loop

generate simulations from markedly different distributions. Algorithm 5 addresses this problem.

A naive approach would be to start iteration t by tuning $d^t(\cdot, \cdot)$ using an additional set of simulations based on parameters drawn from the current importance distribution. However this imposes an additional cost. Instead the algorithm makes a single large set of simulations. These are first used to construct the t th distance function. Then the best N simulations are accepted and used to construct the next importance distribution.

A complication is deciding how many simulations to make for this large set. There must be enough that N of them are accepted. However the distance function defining the acceptance rule is not known until after the simulations are performed. The solution implemented is to continue simulating until $M = \lceil N/\alpha \rceil$ simulations pass the acceptance rule of the *previous* iteration. Let \mathcal{A} be the set of these simulations and \mathcal{B} be the others. Next the new distance function is constructed (based on $\mathcal{A} \cup \mathcal{B}$) and the N with lowest distances (from \mathcal{A}) are accepted. The tuning parameter α has a similar interpretation to the corresponding parameter in Algorithms 3 and 4: the acceptance threshold in iteration t is the α quantile of the realised distances from simulations in \mathcal{A} .

Using this approach means that, as well as adapting the distance function, another difference with Algorithms 3 and 4 is that selection of h_t is delayed from the end of iteration $t - 1$ to part-way through iteration t (and therefore h_1 does not need to be specified as a tuning choice). If desired, this novelty can be used without adapting the distance function. Such a variant of Algorithm 3 was tried on the examples of this paper, but the results are omitted as performance is very similar to Algorithm 3.

Given the same importance density and acceptance rule, an iteration of Algorithm 5 requires the same expected number of simulations as Algorithms 3 and 4. In this sense their costs are the same. In practice, the algorithms select their importance density and acceptance rules differently so this comparison of their computational costs is limited. Section 5 contains empirical comparisons in terms of the mean squared error for a given number of simulations.

Algorithm 5 ABC-PMC with adaptive h_t and $d^t(\cdot, \cdot)$

Initialisation

1. Let $t = 1$.

Main loop

2. Repeat following steps until there are $M = \lceil N/\alpha \rceil$ acceptances.
 - (a) Sample θ^* from importance density $q_t(\theta)$ given in equation (2).
 - (b) If $\pi(\theta^*) = 0$ reject and return to (a).
 - (c) Sample \mathbf{y}^* from $\pi(\mathbf{y}|\theta_i^*)$ and calculate $\mathbf{s}^* = S(\mathbf{y}^*)$.
 - (d) If $t = 1$ accept. Otherwise accept if $d^i(\mathbf{s}^*, \mathbf{s}_{\text{obs}}) \leq h_i$ for all $i < t$.

Denote the accepted parameters as $\theta_1^*, \dots, \theta_M^*$ and the corresponding summary vectors as $\mathbf{s}_1^*, \dots, \mathbf{s}_M^*$.

3. Calculate $(\sigma_1^t, \sigma_2^t, \dots)$, a vector of MADs for each summary statistic, calculated from all the simulations in step 2 (including those rejected).
4. Define $d^t(\cdot, \cdot)$ as the distance (3) using weights $(\omega_i^t)_{1 \leq i \leq m}$ where $\omega_i^t = 1/\sigma_i^t$.
5. Calculate $d_i^* = d^t(\mathbf{s}_i^*, \mathbf{s}_{\text{obs}})$ for $1 \leq i \leq M$.
6. Let h_t be the N th smallest d_i^* value.
7. Let $(\theta_i^t)_{1 \leq i \leq N}$ be the θ_i^* vectors with the smallest d_i^* values (breaking ties randomly).
8. Let $w_i^t = \pi(\theta_i^t)/q_t(\theta_i^t)$ for $1 \leq i \leq N$.
9. Increment t to $t + 1$.

End of loop

The comments at the end of Section 4.1 on tuning details and storing \mathbf{s}^* vectors also apply to Algorithm 5.

4.3 Convergence

This section shows that conditions for the convergence of Algorithms 4 and 5 in practice are essentially those described in Section 2.3 for standard ABC-PMC plus one extra requirement: $e_t = \frac{\max_i w_i^t}{\min_i w_i^t}$ is bounded above.

In more detail, conditions ensuring convergence of Algorithms 4 and 5 can be taken from Theorem 1 in the supplementary material. These are the same as those given for other ABC-PMC algorithms in Section 2.3 with the exception that the acceptance region A_t is now defined as $\{\mathbf{s} | d^t(\mathbf{s}, \mathbf{s}_{\text{obs}}) \leq h_i \text{ for all } i \leq t\}$. Two conditions behave differently under this change: C4 and C5.

Condition C4 states that $\lim_{t \rightarrow \infty} |A_t| = 0$ i.e. Lebesgue measure tends to zero. The definition of A_t for Algorithms 4 and 5 ensures $|A_t|$ is decreasing in t . However it may not converge to zero. Reasons for this are the same as why condition C4 can fail for standard ABC-PMC, as described in Section 2.3.

Condition C5 is bounded eccentricity (defined in the supplementary material) of the A_t sets. Under distance (3) this can easily be seen to correspond to e_t having an upper bound. This is not guaranteed by Algorithms 4 and 5, but it can be imposed, for example by updating ω_i^t to $\omega_i^t + \delta \max_i \omega_i^t$ after step 4 for some small $\delta > 0$. However this was not found to be necessary in any of the examples of this paper.

5 Examples

This section presents three examples comparing the proposed and existing ABC-PMC algorithms: a simple illustrative normal model, the g -and- k distribution and the Lotka–Volterra model.

5.1 Normal distribution

Suppose there is a single parameter θ with prior distribution $N(0, 100^2)$. Let $s_1 \sim N(\theta, 0.1^2)$ and $s_2 \sim N(0, 1^2)$ independently. These are respectively informative and uninformative summary statistics. Let $s_{\text{obs},1} = s_{\text{obs},2} = 0$.

Figures 2 and 3 illustrate the behaviour of ABC-PMC for this example using Algorithms 2 (with adaptive choice of h_t), 4 and 5. For ease of comparison the algorithms use the same random seed, and the distance function and first threshold value h_1 for Algorithms 2 and 4 are specified to be those produced in the first iteration of Algorithm 5. The effect is similar to making a short preliminary run of ABC-rejection to make these tuning choices. All algorithms use $N = 2000$ and $\alpha = 1/2$. (Empirical tests show that $\alpha \approx 1/2$ minimises mean squared error for all algorithms in this and the following examples.)

Under the prior predictive distribution the MAD for s_1 is in the order of 100 while that for s_2 is in the order of 1. Therefore the first acceptance region in Figure 2 is a wide ellipse. Under Algorithm 2 the subsequent acceptance regions are smaller ellipses with the same shape and centre. The acceptance regions for Algorithms 4 and 5 are similar for the first few iterations. After this, enough has been learnt about θ that the simulated summary statistics have a different distribution, with a reduced MAD for s_1 . Hence s_1 is given a larger weight, while the MAD and weight of s_2 remain roughly unchanged. Thus the acceptance regions change shape to become narrower ellipses, which results in a more accurate estimation of θ , as shown by the comparison of mean squared errors (MSEs) in Figure 3. Note that Algorithm 5 adapts its weights more quickly than Algorithm 4 and hence achieves a smaller MSE.

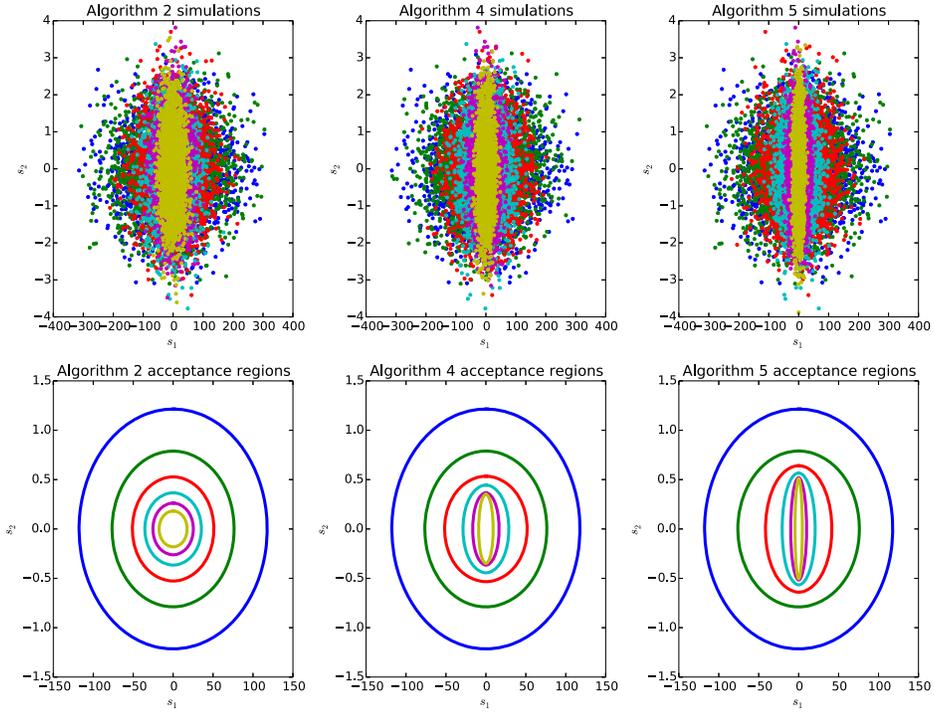


Figure 2: An illustration of ABC-PMC for a simple normal model using Algorithms 2 (non-adaptive distance function), 4 and 5 (adaptive distance functions). *Top row*: simulated summary statistics (including rejections) *Bottom row*: acceptance regions (note different scale to top row). In both rows colour indicates the iteration of the algorithm.

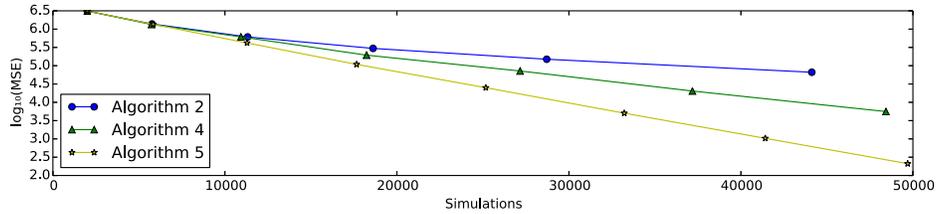


Figure 3: Mean squared error of the parameter for a simple normal example using Algorithms 2, 4 and 5.

5.2 g -and- k distribution

The g -and- k distribution is a popular test of ABC methods. It is defined by its quantile function:

$$A + B \left[1 + c \frac{1 - \exp(-gz(x))}{1 + \exp(-gz(x))} \right] [1 + z(x)^2]^k z(x), \quad (4)$$

where $z(x)$ is the quantile function of the standard normal distribution. Following the literature (Rayner and MacGillivray, 2002), $c = 0.8$ is used throughout. This leaves (A, B, g, k) as unknown parameters.

The g -and- k distribution does not have a closed form density function making likelihood-based inference difficult. However simulation is straightforward: sample $x \sim \text{Unif}(0, 1)$ and substitute into (4). The following example is taken from Drovandi and Pettitt (2011b). Suppose a dataset is 10,000 independent identically distributed draws from the g -and- k distribution and the summary statistics are a subset of the order statistics: those with indices $(1250, 2500, \dots, 8750)$. (As in Fearnhead and Prangle, 2012, a fast method is used to simulate these order statistics without sampling an entire dataset.) The parameters are taken to have independent $\text{Unif}(0, 10)$ priors.

To use as observations, 100 datasets are simulated from the prior predictive distribution. Each is analysed using Algorithms 3, 4 and 5. All analyses uses a total of 10^6 simulations and tuning parameters $N = 1000$ and $\alpha = 1/2$. Table 1 shows root mean squared errors for the output of the algorithms, averaged over all the observed datasets. These show that the adaptive algorithms, 4 and 5, are more accurate overall for every parameter, and perform very similarly to each other.

	A	B	g	k
Algorithm 3	0.335	0.501	0.880	0.163
Algorithm 4	0.083	0.371	0.532	0.126
Algorithm 5	0.081	0.373	0.523	0.126

Table 1: Root mean squared errors of each parameter in the g -and- k example, averaged over analyses of 100 simulated datasets.

More detail is now given for a particular observed dataset, simulated under parameter values $(3, 1, 1.5, 0.5)$. Figure 4 shows the estimated MSE of each parameter for each iteration of the three algorithms. The adaptive algorithms, 4 and 5, performs better throughout for the g and k parameters. For this dataset all the algorithms perform similarly for the location and scale parameters A and B , which have smaller MSE values. Table 2 demonstrates that the main difference in the final estimated posteriors is that Algorithm 3 has higher variances for the g and k parameters.

Figure 5 shows some of the distance function weights produced by the algorithms. Algorithm 3 places low weights on the most extreme order statistics, as they are highly

	A	B	g	k
Algorithm 3	2.98 (0.012)	0.98 (0.028)	1.52 (0.086)	0.50 (0.081)
Algorithm 4	2.98 (0.012)	0.97 (0.025)	1.56 (0.048)	0.53 (0.035)
Algorithm 5	2.98 (0.012)	0.98 (0.024)	1.56 (0.046)	0.53 (0.033)

Table 2: Estimated marginal posterior means and standard deviations (in brackets) of each parameter in the g -and- k example, for analysis of a particular simulated dataset. The values are taken from the final iteration of each algorithm. (n.b. All the estimated posteriors are roughly normal.)

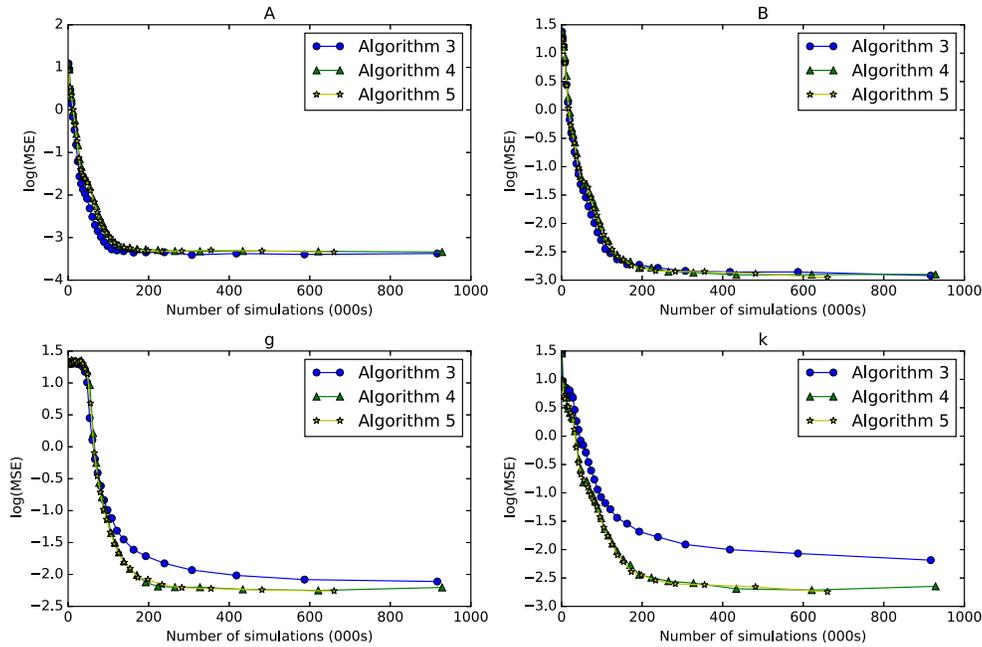
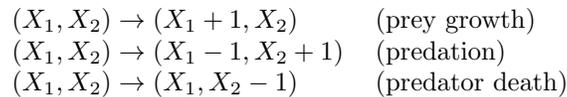


Figure 4: Mean squared error of each parameter from Algorithms 3, 4 and 5 for the g -and- k example.

variable in the prior predictive distribution. This is because the prior places significant weight upon parameter values producing very heavy tails. However by the last iteration of Algorithms 4 and 5 such parameter values have been ruled out. The algorithm therefore assigns larger weights which provide access to the informational content of these statistics.

5.3 Lotka–Volterra model

The Lotka–Volterra model describes two interacting populations. In its original ecological setting the populations represent predators and prey. However it is also a simple example of biochemical reaction dynamics of the kind studied in systems biology. This section concentrates on a stochastic Markov jump process version of this model with state $(X_1, X_2) \in \mathbb{Z}^2$ representing prey and predator population sizes. Three transitions are possible:



These have hazard rates $\theta_1 X_1$, $\theta_2 X_1 X_2$ and $\theta_3 X_2$ respectively. Simulation is straightforward by the Gillespie method. Following either a transition at time t , or initiation at

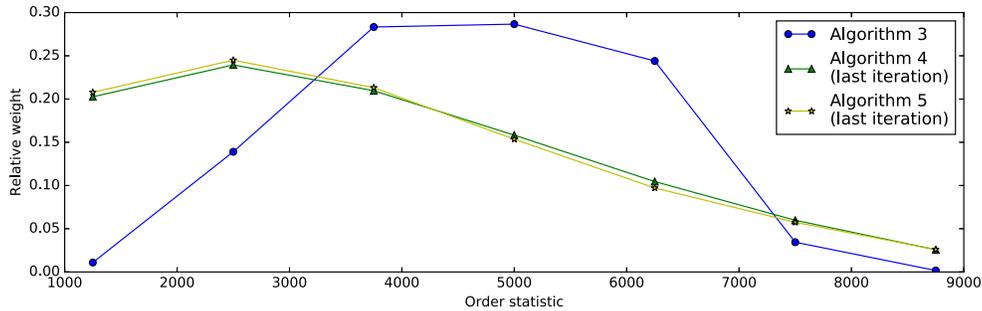


Figure 5: Summary statistic weights used in Algorithms 3, 4 and 5 for the g -and- k example, rescaled to sum to 1.

$t = 0$, the time to the next transition is exponentially distributed with rate equal to the sum of the hazard rates at time t . The type of the next transition has a multinomial distribution with probabilities proportional to the hazard rates. For more background see for example Owen et al. (2015), from which the following specific inference problem is taken.

The initial conditions are taken to be $X_1 = 50, X_2 = 100$. A dataset is formed of observations at times $2, 4, 6, \dots, 32$. Both X_1 and X_2 are observed plus independent $N(0, \sigma^2)$ errors, where σ is fixed at $\exp(2.3)$. The unknown parameters are taken to be $\log \theta_1, \log \theta_2$ and $\log \theta_3$. These are given independent $\text{Unif}(-6, 2)$ priors. The vector of all 32 noisy observations is used as the ABC summary statistics.

A single simulated dataset is analysed (shown in Figure 8). This is generated from the model with $\theta_1 = 1, \theta_2 = 0.005, \theta_3 = 0.6$. ABC analysis is performed using Algorithms 3, 4 and 5. A total of 50,000 simulations are used in each algorithm. The tuning parameters are $N = 200$ and $\alpha = 1/2$. Any Lotka–Volterra simulation reaching 100,000 transitions is terminated and automatically rejected. This avoids extremely long simulations, such as exponential prey growth if predators die out. These incomplete simulations are excluded from the MAD calculations, but this should have little effect as they are rare.

Figure 6 shows the MSEs resulting from the analyses. The adaptive algorithms, 4 and 5, have similar outputs. Both produce smaller errors than Algorithm 3 for all parameters after roughly 10,000 simulations. Table 3 demonstrates that the main difference in the final estimated posteriors is that Algorithm 3 has higher variances. Figure 7 shows the weights used throughout Algorithm 3 and those used in the final iteration of the others. Again the adaptive algorithms are similar to each other but different to Algorithm 3. Figure 8 explains this by showing a sample of simulated datasets on which these weights are based. Under the prior predictive distribution (shown in the top row), at least one population usually quickly becomes extinct, illustrating that the prior distribution concentrates on the wrong system dynamics and so is unsuitable for choosing distance weights for later iterations of the algorithm.

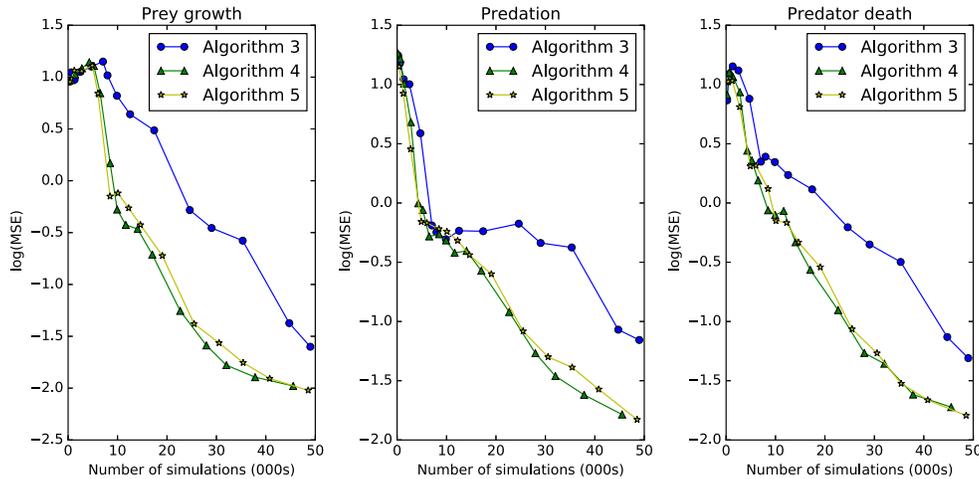


Figure 6: Mean squared error of each parameter (i.e. $\log \theta_1, \log \theta_2, \log \theta_3$) from ABC-PMC output for the Lotka–Volterra example.

	$\log \theta_1$	$\log \theta_2$	$\log \theta_3$
Algorithm 3	-0.048 (0.15)	-5.15 (0.21)	-0.48 (0.22)
Algorithm 4	-0.021 (0.10)	-5.24 (0.11)	-0.56 (0.13)
Algorithm 5	-0.021 (0.10)	-5.24 (0.11)	-0.55 (0.12)

Table 3: Estimated marginal posterior means and standard deviations (in brackets) of each parameter in the Lotka–Volterra example, for analysis of a particular simulated dataset. The values are taken from the final iteration of each algorithm. The true values are 0, -5.30 and -0.51 . (n.b. All the estimated posteriors are roughly normal.)

6 Discussion

This paper has presented two ABC-PMC algorithms with adaptive distance functions. The algorithms adapt the structure to ABC-PMC by using the output of existing simulation steps to adapt their distance functions. Therefore they have a similar computational cost for the same number of iterations. Furthermore, their convergence properties are similar to ABC-PMC. Several examples have been shown where the new algorithms improve performance. This is because in each example the scale of the summary statistics varies significantly between prior and posterior predictive distributions. Of the two algorithms, Algorithm 4 is simpler to implement, involving only a small modification to standard ABC-PMC, and has essentially the same performance to Algorithm 5 in two of the three examples. Algorithm 5 performs better in the example of Section 5.1, suggesting it is preferable in situation where continual adaptation is required. The remainder of this section discusses possibilities to extend this work.

Several variations on ABC-PMC have been proposed in the literature. The adaptive distance function idea introduced here can be used in most of these. This is particu-

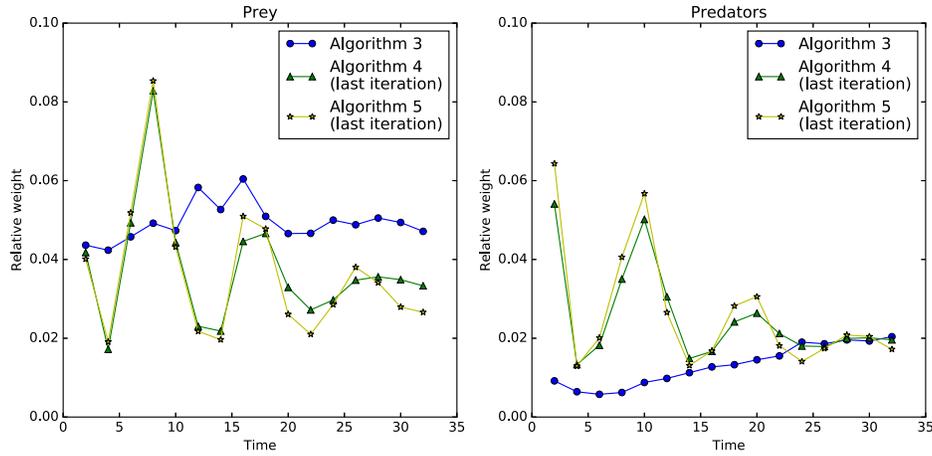


Figure 7: Summary statistic weights used in ABC-PMC for the Lotka–Volterra example, rescaled to sum to 1.

larly simple for ABC model choice algorithms (e.g. Toni et al., 2009). Here, instead of proposing θ^* values from an importance density, (m^*, θ^*) pairs are proposed, where m^* is a model indicator. This could be implemented in Algorithms 4 and 5 while leaving the other details unchanged. Drovandi and Pettitt (2011a), Del Moral et al. (2012) and Lenormand et al. (2013) propose ABC-SMC algorithms which update the population of (θ, s) pairs between iterations in different ways to ABC-PMC. In all of these it seems possible to update distance functions using the strategies of Algorithms 4 and 5. However some of these variations would require further convergence results beyond those given in the supplementary material.

Several aspects of Algorithms 4 and 5 could be modified. One natural alternative is to use Mahalanobis-style distance functions $d^t(\mathbf{x}, \mathbf{y}) = [(\mathbf{x} - \mathbf{y})^T W^t (\mathbf{x} - \mathbf{y})]^{1/2}$ where W^t is an estimate of the precision matrix. Scenarios exist in which this performs much better than weighted Euclidean distance, (3) (Sisson, personal communication). However exploratory work found it gave similar or worse performance for the examples in this paper. Distance (3) is preferred here for this reason, and also because its weights are easier to interpret and there are more potential numerical difficulties in estimating a precision matrix. Nonetheless, for other problems it may be worth considering both alternatives.

Another reason it may be desirable to modify the distance function (3) is if some summary statistic, say s_i , has an observed value far from most simulated values. In this case $|s_{\text{obs},i} - s_i|$ can be much larger than σ_i , and so s_i can dominate the distances used in this paper. It is tempting to downweight s_i so that the others summaries can also contribute. Finding a good way to do this without ignoring s_i altogether is left for future work.

Algorithms 4 and 5 update the distance function at each iteration. There may be scope for similarly updating other tuning choices. It is particularly appealing to try

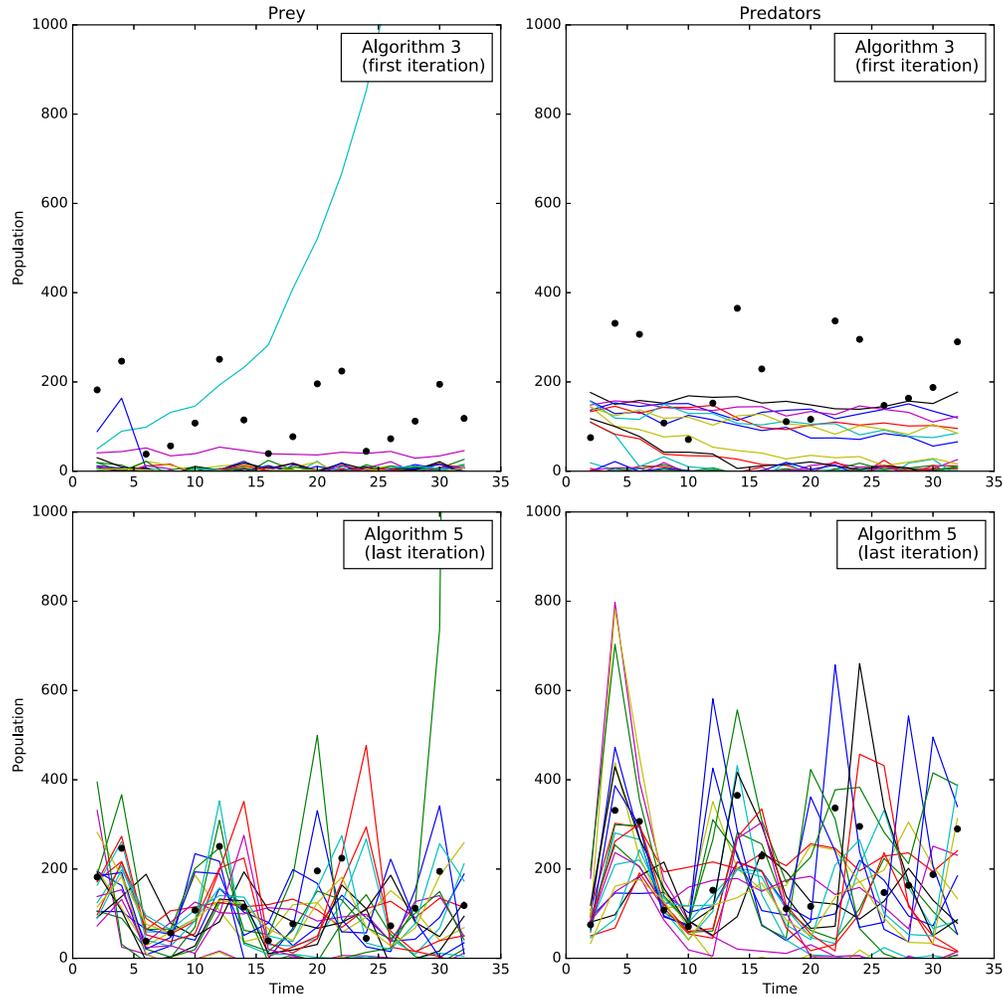


Figure 8: Observed dataset (black points) and samples of 20 simulated datasets (coloured lines) for the Lotka–Volterra example. The top row shows simulations from step 2 of the first iteration of Algorithm 3. The bottom row shows simulations from step 2 of the last iteration of Algorithm 5. These are representative examples of the simulations used to select the weights shown in Figure 7. Simulations for Algorithm 4 are not shown but are qualitatively similar to the bottom row.

to improve the choice of summary statistics as the algorithm progresses (as suggested by Barnes et al., 2012). Summary statistics could be selected at the same time as the distance function based on the same simulations, for example by a modification of the regression method of Fearnhead and Prangle (2012). Further work would be required to ensure the convergence of such an algorithm.

Supplementary Material

Adapting the ABC distance function: Supplementary Material
(DOI: [10.1214/16-BA1002SUPP](https://doi.org/10.1214/16-BA1002SUPP); .pdf).

References

- Barnes, C. P., Filippi, S., and Stumpf, M. P. H. (2012). “Contribution to the discussion of Fearnhead and Prangle (2012).” *Journal of the Royal Statistical Society: Series B*, 74: 453. [306](#)
- Beaumont, M. A. (2010). “Approximate Bayesian computation in evolution and ecology.” *Annual Review of Ecology, Evolution and Systematics*, 41: 379–406. [291](#)
- Beaumont, M. A., Cornuet, J.-M., Marin, J.-M., and Robert, C. P. (2009). “Adaptive approximate Bayesian computation.” *Biometrika*, 2025–2035. [MR2767283](#). doi: <http://dx.doi.org/10.1093/biomet/asp052>. [290](#), [291](#)
- Beaumont, M. A., Zhang, W., and Balding, D. J. (2002). “Approximate Bayesian Computation in Population Genetics.” *Genetics*, 162: 2025–2035. [289](#), [294](#)
- Bezanson, J., Karpinski, S., Shah, V. B., and Edelman, A. (2012). “Julia: A fast dynamic language for technical computing.” [arXiv:1209.5145](#). [290](#)
- Biau, G., Cérou, F., and Guyader, A. (2015). “New insights into Approximate Bayesian Computation.” *Annales de l’Institut Henri Poincaré (B) Probabilités et Statistiques*, 51(1): 376–403. [MR3300975](#). doi: <http://dx.doi.org/10.1214/13-AIHP590>. [291](#)
- Blum, M. G. B., Nunes, M. A., Prangle, D., and Sisson, S. A. (2013). “A comparative review of dimension reduction methods in approximate Bayesian computation.” *Statistical Science*, 28: 189–208. [MR3112405](#). [295](#)
- Bonassi, F. V. and West, M. (2015). “Sequential Monte Carlo with Adaptive Weights for Approximate Bayesian Computation.” *Bayesian Analysis*, 10(1): 171–187. [MR3420901](#). doi: <http://dx.doi.org/10.1214/14-BA891>. [290](#)
- Cappé, O., Guillin, A., Marin, J.-M., and Robert, C. P. (2004). “Population Monte Carlo.” *Journal of Computational and Graphical Statistics*, 13(4). [MR2109057](#). doi: <http://dx.doi.org/10.1198/106186004X12803>. [291](#)
- Csilléry, K., Blum, M. G. B., Gaggiotti, O., and François, O. (2010). “Approximate Bayesian Computation in practice.” *Trends in Ecology & Evolution*, 25: 410–418. [291](#)
- Csilléry, K., François, O., and Blum, M. G. B. (2012). “abc: an R package for approximate Bayesian computation (ABC).” *Methods in Ecology and Evolution*, 3: 475–479. [294](#)
- Del Moral, P., Doucet, A., and Jasra, A. (2012). “An adaptive sequential Monte Carlo method for approximate Bayesian computation.” *Statistics and Computing*, 22(5): 1009–1020. [MR2950081](#). doi: <http://dx.doi.org/10.1007/s11222-011-9271-y>. [292](#), [305](#)

- Drovandi, C. C. and Pettitt, A. N. (2011a). “Estimation of parameters for macroparasite population evolution using approximate Bayesian computation.” *Biometrics*, 67(1): 225–233. MR2898834. doi: <http://dx.doi.org/10.1111/j.1541-0420.2010.01410.x>. 292, 305
- Drovandi, C. C. and Pettitt, A. N. (2011b). “Likelihood-free Bayesian estimation of multivariate quantile distributions.” *Computational Statistics & Data Analysis*, 55(9): 2541–2556. MR2802334. doi: <http://dx.doi.org/10.1016/j.csda.2011.03.019>. 301
- Fasiolo, M. and Wood, S. N. (2015). “Approximate methods for dynamic ecological models.” [arXiv:1511.02644](https://arxiv.org/abs/1511.02644). 290
- Fearnhead, P. and Prangle, D. (2012). “Constructing Summary Statistics for Approximate Bayesian Computation: Semi-automatic ABC.” *Journal of the Royal Statistical Society, Series B*, 74: 419–474. MR2925370. doi: <http://dx.doi.org/10.1111/j.1467-9868.2011.01010.x>. 301, 306
- Lenormand, M., Jabot, F., and Deffuant, G. (2013). “Adaptive approximate Bayesian computation for complex models.” *Computational Statistics*, 28(6): 2777–2796. MR3141363. doi: <http://dx.doi.org/10.1007/s00180-013-0428-3>. 305
- Marin, J.-M., Pudlo, P., Robert, C. P., and Ryder, R. J. (2012). “Approximate Bayesian computational methods.” *Statistics and Computing*, 22(6): 1167–1180. MR2992292. doi: <http://dx.doi.org/10.1007/s11222-011-9288-2>. 291
- McKinley, T., Cook, A. R., and Deardon, R. (2009). “Inference in epidemic models without likelihoods.” *The International Journal of Biostatistics*, 5(1). MR2533810. doi: <http://dx.doi.org/10.2202/1557-4679.1171>. 294
- Owen, J., Wilkinson, D. J., and Gillespie, C. S. (2015). “Likelihood free inference for Markov processes: a comparison.” *Statistical applications in genetics and molecular biology*, 14(2): 189–209. MR3331773. doi: <http://dx.doi.org/10.1515/sagmb-2014-0072>. 294, 303
- Prangle, D. (2016). “Adapting the ABC distance function: Supplementary Material.” *Bayesian Analysis*. doi: <http://dx.doi.org/10.1214/16-BA1002SUPP>. 290
- Rayner, G. D. and MacGillivray, H. L. (2002). “Numerical maximum likelihood estimation for the g-and-k and generalized g-and-h distributions.” *Statistics and Computing*, 12(1): 57–75. MR1877580. doi: <http://dx.doi.org/10.1023/A:1013120305780>. 301
- Sedki, M., Pudlo, P., Marin, J.-M., Robert, C. P., and Cornuet, J.-M. (2012). “Efficient learning in ABC algorithms.” [arXiv:1210.1388](https://arxiv.org/abs/1210.1388). 290
- Silk, D., Filippi, S., and Stumpf, M. P. H. (2013). “Optimizing threshold-schedules for sequential approximate Bayesian computation: applications to molecular systems.” *Statistical Applications in Genetics and Molecular Biology*, 12(5): 603–618. MR3108049. doi: <http://dx.doi.org/10.1515/sagmb-2012-0043>. 292

- Sisson, S. A., Fan, Y., and Tanaka, M. M. (2009). “Correction: Sequential Monte Carlo without likelihoods.” *Proceedings of the National Academy of Sciences*, 106(39): 16889–16890. [290](#), [291](#)
- Toni, T., Welch, D., Strelkowa, N., Ipsen, A., and Stumpf, M. (2009). “Approximate Bayesian computation scheme for parameter inference and model selection in dynamical systems.” *Journal of The Royal Society Interface*, 6(31): 187–202. [290](#), [291](#), [305](#)