



HAL
open science

Approximate Bayesian computation with surrogate posteriors

Florence Forbes, Hien Duy Nguyen, Trung Tin Nguyen, Julyan Arbel

► **To cite this version:**

Florence Forbes, Hien Duy Nguyen, Trung Tin Nguyen, Julyan Arbel. Approximate Bayesian computation with surrogate posteriors. 2021. hal-03139256v3

HAL Id: hal-03139256

<https://hal.science/hal-03139256v3>

Preprint submitted on 26 Aug 2021 (v3), last revised 25 Sep 2022 (v5)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Approximate Bayesian computation with surrogate posteriors

Florence Forbes^{1*}, Hien Duy Nguyen², TrungTin Nguyen³ and Julyan Arbel¹

¹Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK, Inria Grenoble
Rhone-Alpes, 655 av. de l'Europe, 38335 Montbonnot, France.

²School of Engineering and Mathematical Sciences, La Trobe University, Bundoora,
Victoria, Australia.

³Normandie Univ, UNICAEN, CNRS, LMNO, Caen, 14000, France .

*Corresponding author(s). E-mail(s): florence.forbes@inria.fr;

Contributing authors: H.Nguyen5@latrobe.edu.au; trung-tin.nguyen@unicaen.fr;
julyan.arbel@inria.fr;

Abstract

A key ingredient in approximate Bayesian computation (ABC) procedures is the choice of a discrepancy that describes how different the simulated and observed data are, often based on a set of summary statistics when the data cannot be compared directly. Unless discrepancies and summaries are available from experts or prior knowledge, which seldom occurs, they have to be chosen and this can affect the quality of approximations. The choice between discrepancies is an active research topic, which has mainly considered data discrepancies requiring samples of observations or distances between summary statistics. In this work, we introduce a preliminary learning step in which surrogate posteriors are built from finite Gaussian mixtures, using an inverse regression approach. These surrogate posteriors are then used in place of summary statistics and compared using metrics between distributions in place of data discrepancies. Two such metrics are investigated, a standard L_2 distance and an optimal transport-based distance. The whole procedure can be seen as an extension of the semi-automatic ABC framework to functional summary statistics setting and can also be used as an alternative to sample-based approaches. The resulting ABC quasi-posterior distribution is shown to converge to the true one, under standard conditions. Performance is illustrated on both synthetic and real data sets, where it is shown that our approach is particularly useful, when the posterior is multimodal.

Keywords: Approximate Bayesian computation, summary statistics, surrogate models, Gaussian mixtures, Wasserstein distance, multimodal posterior distributions.

1 Introduction

Approximate Bayesian computation (ABC) (see, *e.g.*, Sisson et al. 2019) appears as a natural candidate for addressing problems, where there is a lack of availability or tractability of the likelihood. Such cases occur when the direct model or data generating process is not available analytically, but is available as a simulation procedure; *e.g.*, when the data generating process is characterized as a series of ordinary differential equations, as in Mesejo et al. (2016); Hovorka et al. (2004). In addition, typical features or constraints that can occur in practice are that: 1) the observations \mathbf{y} are high-dimensional, because they represent signals in time or spectra, as in Schmidt and Fernando (2015); Bernard-Michel et al. (2009); Ma et al. (2013); and 2) the parameter $\boldsymbol{\theta}$, to be estimated, is itself multi-dimensional with correlated dimensions so that independently predicting its components is sub-optimal; *e.g.*, when there are known constraints such as when the parameter elements are concentrations or probabilities that sum to one (Deleforge et al., 2015; Lemasson et al., 2016; Bernard-Michel et al., 2009).

The fundamental idea of ABC is to generate parameter proposals $\boldsymbol{\theta}$ in a parameter space Θ using a prior distribution $\pi(\boldsymbol{\theta})$ and accept a proposal if the simulated data \mathbf{z} for that proposal is similar to the observed data \mathbf{y} , both in an observation space \mathcal{Y} . This similarity is usually measured using a distance or discriminative measure D and a simulated sample \mathbf{z} is retained if $D(\mathbf{z}, \mathbf{y})$ is smaller than a given threshold ϵ . In this simple form, the procedure is generally referred to as rejection ABC. Other variants are possible and often recommended, for instance using MCMC or sequential procedures (*e.g.*, Del Moral et al., 2012; Buchholz and Chopin, 2019). We will focus on the rejection version for the purpose of this paper as all developments can be easily adapted to more sophisticated variants. Also to our knowledge, theoretical results only exist for rejection ABC.

In the case of a rejection algorithm, selected samples are drawn from the so-called ABC quasi-posterior, which is an approximation to the true posterior $\pi(\boldsymbol{\theta} \mid \mathbf{y})$. Under conditions similar to that of Bernton et al. (2019), regarding the existence of a probability density function (pdf) $f_{\boldsymbol{\theta}}(\mathbf{z})$ for the likelihood, the ABC quasi-posterior depends on D and on a threshold ϵ , and can be written as

$$\pi_{\epsilon}(\boldsymbol{\theta} \mid \mathbf{y}) \propto \pi(\boldsymbol{\theta}) \int_{\mathcal{Y}} \mathbf{1}_{\{D(\mathbf{y}, \mathbf{z}) \leq \epsilon\}} f_{\boldsymbol{\theta}}(\mathbf{z}) d\mathbf{z} . \quad (1)$$

More specifically, the similarity between \mathbf{z} and \mathbf{y} is generally evaluated based on two components: the choice of summary statistics $s(\cdot)$ to account for the data in a more robust manner, and the choice of a distance to compare the summary statistics. That is, $D(\mathbf{y}, \mathbf{z})$ in (1) should then be replaced by $D(s(\mathbf{y}), s(\mathbf{z}))$, whereupon we overload D to also denote the distance between summary statistics $s(\cdot)$.

However, there is no general rule for constructing good summary statistics for complex models and if a summary statistic does not capture important characteristics of the data, the ABC algorithm is likely to yield samples from an incorrect posterior (Blum et al., 2013; Fearnhead and Prangle, 2012; Gutmann et al., 2018). Great insight has been gained through the work of Fearnhead and Prangle (2012), who introduced the *semi-automatic* ABC framework and showed that under a quadratic loss, the optimal choice for the summary statistic of \mathbf{y} was the true posterior mean of the parameter: $s(\mathbf{y}) = \mathbb{E}[\boldsymbol{\theta} \mid \mathbf{y}]$. This conditional expectation cannot be calculated analytically but can be estimated by regression using a learning data set prior to the ABC procedure itself.

In Fearnhead and Prangle (2012), it is suggested that a simple regression model may be enough to approximate $\mathbb{E}[\boldsymbol{\theta} \mid \mathbf{y}]$, but this has since been contradicted, for instance by Jiang et al. (2017) and Wqvist et al. (2019), who show that the quality of the approximation can matter in practice. Still

focusing on posterior means as summary statistics, they use deep neural networks that capture complex non-linear relationships and exhibit much better results than standard regression approaches. However, deep neural networks remain very computationally costly tools, both in terms of the required size of training data and number of parameters and hyperparameters to be estimated and tuned.

Our first contribution is to investigate an alternative efficient way to construct summary statistics, in the same vein as semi-automatic ABC, but based on posterior moments, not restricted to the posterior means. Although this natural extension was already proposed in Jiang et al. (2017), it requires the availability of a flexible and tractable regression model, able to capture complex non-linear relationships and to provide posterior moments, straightforwardly. As such, Jiang et al. (2017) did not consider an implementation of the procedure. For this purpose, the Gaussian Locally Linear Mapping (GLLiM) method (Deleforge et al., 2015), that we recall in Section 3, appears as a good candidate, with properties that balance between the computationally expensive neural networks and the simple standard regression techniques. In contrast to most regression methods that provide only pointwise predictions, GLLiM provides, at low cost, a parametric estimation of the full true posterior distributions. Using a learning set of parameters and observations couples, GLLiM learns a family of finite Gaussian mixtures whose parameters depend analytically on the observation to be inverted. For any observed data, the true posterior can be approximated as a Gaussian mixture, whose moments are easily computed and turned into summary statistics for subsequent ABC sample selection.

Our second contribution is to propose to compare directly the full surrogate posterior distributions provided by GLLiM, without reducing them to their moments. So doing, we introduce the idea of functional summary statistics, which also requires a different notion of the usual distances or discrepancy measures to compare them. Recent developments in optimal transport-based distances designed for Gaussian mixtures (Delon and Desolneux, 2020; Chen et al., 2019) match perfectly this need via the so-called Mixture-Wasserstein distance as referred to by Delon and Desolneux (2020), and denoted throughout the text as MW_2 . There exist other distances between mixtures that are tractable, and among them the L_2 distance is also considered in this work.

A remarkable feature of our approach is that it can be equally applied to settings where a sample of *i.i.d.* observations is available (*e.g.* Bernton et al. (2019); Nguyen et al. (2020)) and to settings where a single observation is available, as a vector of measures, a time series realization or a data set reduced to a vector of summary statistics (*e.g.* Fearnhead and Prangle (2012); Drovandi and Pettitt (2011)).

The novelty of our approach and its comparison with existing work is emphasized in Section 2. The GLLiM output is briefly described in Section 3. A first exploitation of GLLiM combined with the semi-automatic ABC principle is presented in Section 4.1. Our extension, using functional summary statistics, is then described in Section 4.2. The approach’s theoretical properties are investigated in Section 5 and the practical performance is illustrated in Section 6, both on synthetic and real data. Detailed proofs and additional illustrations are shown in a supplementary material file. The code can be found at <https://github.com/Trung-TinNguyenDS/GLLiM-ABC>.

2 Related work

As an alternative to semi-automatic ABC, in the works of Nguyen et al. (2020); Jiang et al. (2018); Bernton et al. (2019); Park et al. (2016); Gutmann et al. (2018), the difficulties associated with finding efficient summary statistics were bypassed by adopting, respectively, the Energy Distance, a

Kullback–Leibler divergence estimator, the Wasserstein distance, the Maximum Mean Discrepancy (MMD), and classification accuracy to provide a data discrepancy measure. Such approaches compare simulated data and observed data by looking at them as *i.i.d.* samples from distributions, respectively linked to the simulated and true parameter, except for Bernton et al. (2019) and Gutmann et al. (2018) who proposed solutions to also handle time series. These methods require sufficiently large samples and cannot be applied if the sample related to the parameter to be recovered is too small. This is a major difference with the approach we propose, which can be applied in both cases. We refer to these two cases as *one observation* and *i.i.d. observations* settings. In the *one observation* case, the observed data restricts to a single observation \mathbf{y} of dimension d assumed to be generated from a true parameter $\boldsymbol{\theta}$ of dimension ℓ . This case is commonly encountered in inverse problems where it may be impossible to gather repeated observations from the same parameter values due to technological reasons. Typically, in remote sensing applications, satellites are limited to only a few degrees of freedom when observing a given site in constant conditions. This is also the case when the observation is a time series or when a sample of observations is reduced to a single vector of summary statistics. In the multiple *i.i.d.* observations case, the observed data is made of a sample of R *i.i.d.* realizations $\{\mathbf{y}^1, \dots, \mathbf{y}^R\}$ coming from the same true $\boldsymbol{\theta}$. The previous case is trivially recovered when $R = 1$.

ABC procedures using a regression step, as introduced by Fearnhead and Prangle (2012), are adapted to one observation settings. They cannot be applied on large (*e.g.* $R = 10^4$) numbers of covariates and require that samples, observed and simulated, are first reduced to a smaller number of statistics, *e.g.* 100. In contrast, discrepancy-based approaches compare empirical distributions constructed from the samples and require a relatively large R .

Our method is not limited to either one of these cases because we do not compare samples from distributions, but directly the distributions through their surrogates using distances between distributions. We can use the same Wasserstein, Kullback–Leibler divergence, *etc.*, but in their *population* versions rather than in their empirical versions. A Wasserstein-based distance can be computed between mixtures of Gaussians, thanks to the recent work of Delon and Desolneux (2020) and Chen et al. (2019). Closed form expressions also exist for the L_2 distance, for the MMD with a Gaussian RBF kernel, or a polynomial kernel (see Sriperumbudur et al., 2010; Muandet et al., 2012) and for the Jensen–Rényi divergence of degree two (see Wang et al., 2009). Kristan et al. (2011) also proposed an algorithm based on the so-called insceted transform in order to compute the Hellinger distance between two Gaussian mixtures, although it is unclear what the complexity of this algorithm is.

To emphasize the difference to more standard summaries, we refer to our surrogate posteriors as functional summary statistics. The term has already been used by Soubeyrand et al. (2013) in the ABC context in their attempts to characterize spatial structures using statistics that are functions (*e.g.* correlograms or variograms). They do not address the issue of choosing the summary statistics. Given such functional statistics whose nature may change for each considered model, their goal is to optimize the distances to compare them. In our proposal, the functional statistics are probability distributions. They arise as a way to bypass the summary statistics choice, but in this work, we make use of existing metrics to compare them, without optimization.

3 Parametric posterior approximation with Gaussian mixtures

A learning set $\mathcal{D}_N = \{(\boldsymbol{\theta}_n, \mathbf{y}_n), n \in [N]\}$ is built from the joint distribution that results from the prior $\pi(\boldsymbol{\theta})$ on $\boldsymbol{\theta}$ and the likelihood $f_{\boldsymbol{\theta}}$, where $[N] = \{1, \dots, N\}$. The idea is to capture the relationship

between $\boldsymbol{\theta}$ and \mathbf{y} with a joint probabilistic model for which computing conditional distributions and moments is straightforward. For the choice of the model to fit to \mathcal{D}_N , we propose to use the so-called Gaussian Locally Linear Mapping (GLLiM) model (Deleforge et al., 2015) for its ability to capture non-linear relationships in a tractable manner, based on flexible mixtures of Gaussian distributions. GLLiM can be considered within the class of inverse regression approaches, such as sliced inverse regression (Li, 1991), partial least squares (Cook and Forzani, 2019), mixtures of regressions approaches of different variants, *e.g.* mixtures of experts (Nguyen et al., 2019), cluster weighted models (Ingrassia et al., 2012), and kernel methods (Nataraj et al., 2018). In contrast to deep learning approaches (see Arridge et al. 2019, for a survey), GLLiM provides for each observed \mathbf{y} , a full posterior probability distribution within a family of parametric models $\{p_G(\boldsymbol{\theta} | \mathbf{y}; \boldsymbol{\phi}), \boldsymbol{\phi} \in \Phi\}$. To model non-linear relationships, it uses a mixture of K linear models. More specifically, the expression of $p_G(\boldsymbol{\theta} | \mathbf{y}; \boldsymbol{\phi})$ is analytical and available for all \mathbf{y} with $\boldsymbol{\phi}$ being independent of \mathbf{y} :

$$p_G(\boldsymbol{\theta} | \mathbf{y}; \boldsymbol{\phi}) = \sum_{k=1}^K \eta_k(\mathbf{y}) \mathcal{N}(\boldsymbol{\theta}; \mathbf{A}_k \mathbf{y} + \mathbf{b}_k, \boldsymbol{\Sigma}_k), \quad (2)$$

where $\mathcal{N}(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes the Gaussian pdf with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$ and $\eta_k(\mathbf{y}) = \pi_k \mathcal{N}(\mathbf{y}; \mathbf{c}_k, \boldsymbol{\Gamma}_k) / \sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{y}; \mathbf{c}_j, \boldsymbol{\Gamma}_j)$. This distribution involves parameters: $\boldsymbol{\phi} = \{\pi_k, \mathbf{c}_k, \boldsymbol{\Gamma}_k, \mathbf{A}_k, \mathbf{b}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$. One interesting property of this model is that the mixture setting provides guarantees that, when choosing K large enough, it is possible to approximate any reasonable relationship (Nguyen et al., 2019, 2020). The parameter $\boldsymbol{\phi}$ can be estimated by fitting a GLLiM model to \mathcal{D}_N using an Expectation-Maximization (EM) algorithm. Details are provided in supplementary material and in Deleforge et al. (2015).

Fitting a GLLiM model to \mathcal{D}_N therefore results in a set of parametric distributions $\{p_G(\boldsymbol{\theta} | \mathbf{y}; \boldsymbol{\phi}_{K,N}^*), \mathbf{y} \in \mathcal{Y}\}$, which are mixtures of Gaussian distributions and can be seen as a parametric mapping from \mathbf{y} values to posterior pdfs on $\boldsymbol{\theta}$. The parameter $\boldsymbol{\phi}_{K,N}^*$ is the same for all conditional distributions and does not need to be re-estimated for each new instance of \mathbf{y} . When required, it is straightforward to compute the expectation and covariance matrix of $p_G(\boldsymbol{\theta} | \mathbf{y}; \boldsymbol{\phi}_{K,N}^*)$ in (2):

$$\mathbb{E}_G[\boldsymbol{\theta} | \mathbf{y}; \boldsymbol{\phi}_{K,N}^*] = \sum_{k=1}^K \eta_k^*(\mathbf{y}) (\mathbf{A}_k^* \mathbf{y} + \mathbf{b}_k^*), \quad (3)$$

$$\begin{aligned} \text{Var}_G[\boldsymbol{\theta} | \mathbf{y}; \boldsymbol{\phi}_{K,N}^*] &= \sum_{k=1}^K \eta_k^*(\mathbf{y}) [\boldsymbol{\Sigma}_k^* + (\mathbf{A}_k^* \mathbf{y} + \mathbf{b}_k^*)(\mathbf{A}_k^* \mathbf{y} + \mathbf{b}_k^*)^\top] \\ &\quad - \mathbb{E}_G[\boldsymbol{\theta} | \mathbf{y}; \boldsymbol{\phi}_{K,N}^*] \mathbb{E}_G[\boldsymbol{\theta} | \mathbf{y}; \boldsymbol{\phi}_{K,N}^*]^\top. \end{aligned} \quad (4)$$

Expression (3) then provides approximate posterior means and can be directly used in a semi-automatic ABC procedure. In addition, summary statistics extracted from the covariance matrix (4) can also be included and is likely to improve the ABC selection as illustrated in Section 6.

When R *i.i.d.* d -dimensional observations are available for each parameter value, they can be stacked into a single large vector. However, as noted by Fearnhead and Prangle (2012) and Jiang et al. (2017), the resulting number of covariates, at least $d \times R$, may become too large. Even if this is computationally doable with the standard GLLiM procedure, it is likely to be sub-optimal as it ignores the *i.i.d.* nature of the data. To handle this case, we therefore propose an adaptation of the

EM algorithm of Deleforge et al. (2015). This adaptation is detailed in the supplementary material Section S1 and illustrated in the first two examples of Section 6. It is shown by Deleforge et al. (2015) that constraints on the model parameterization can be assumed without oversimplifying mixture (2). These constraints concern the covariance matrices used in the mixture modeling of the likelihood (or the direct model) and are not directly visible on the Σ_k 's which remain full in general. In addition to model the *i.i.d.* case, the adaptation we propose adds to the existing constraints, isotropic or diagonal matrices, the possibility to assume bloc diagonal structures.

4 Extended semi-automatic ABC

Semi-automatic ABC refers to an approach introduced in Fearnhead and Prangle (2012), which has since then led to various attempts and improvements, see *e.g.* Jiang et al. (2017) and Wiquist et al. (2019), without dramatic deviation from the original ideas.

4.1 Extension to extra summary vectors

A natural idea is to use the approximate posterior expectation provided by GLLiM in (3) as the summary statistic s of data \mathbf{y} , $s(\mathbf{y}) = \mathbb{E}_G[\boldsymbol{\theta} \mid \mathbf{y}; \boldsymbol{\phi}_{K,N}^*]$. It provides a first attempt to combine GLLiM and ABC procedures and has the advantage over neural networks of being easier to estimate without the need for huge learning data sets and hyperparameter tuning.

However, one advantage of GLLiM over most regression methods is not to reduce to pointwise predictions and to provide full posteriors as output. The posteriors can then be used to provide other posterior moments as summary statistics. The same standard ABC procedure as before can be applied but now with $s_1(\mathbf{y}) = \mathbb{E}_G[\boldsymbol{\theta} \mid \mathbf{y}; \boldsymbol{\phi}_{K,N}^*]$ and $s_2(\mathbf{y}) = \text{Var}_G[\boldsymbol{\theta} \mid \mathbf{y}; \boldsymbol{\phi}_{K,N}^*]$, as given by (4).

As illustrated in Section 6, it is easy to construct examples where the posterior expectations, even when well-approximated, do not perform well as summary statistics. Providing a straightforward and tractable way to add other posterior moments is then already an interesting contribution. However, to really make the most of the GLLiM framework, we propose to further exploit the fact that GLLiM provides more than moments.

4.2 Extension to functional summary statistics

Instead of comparing simulated \mathbf{z} 's to the observed \mathbf{y} , or equivalently their summary statistics, we propose to compare the $p_G(\boldsymbol{\theta} \mid \mathbf{z}; \boldsymbol{\phi}_{K,N}^*)$'s to $p_G(\boldsymbol{\theta} \mid \mathbf{y}; \boldsymbol{\phi}_{K,N}^*)$, as given by (2). As approximations of the true posteriors, these quantities are likely to capture the main characteristics of $\boldsymbol{\theta}$ without committing to the choice of a particular moment. The comparison requires an appropriate distance that needs to be a mathematical distance between distributions. The equivalent functional distance to the L_2 distance can still be used, as can the Hellinger distance or any other divergence. A natural choice is the Kullback–Leibler divergence, but computing it between mixtures is not straightforward. Computing the Energy statistic (*e.g.*, Nguyen et al., 2020) appears at first to be easier but in the end that would still resort to Monte Carlo sums. Since model (2) is parametric, we could also compute distances between the parameters of the mixtures that depend on \mathbf{y} . That is for $k \in [K]$, between the mixing proportions $\eta_k^*(\mathbf{y}) = \frac{\pi_k^* \mathcal{N}(\mathbf{y}; \mathbf{c}_k^*, \boldsymbol{\Gamma}_k^*)}{\sum_{j=1}^K \pi_j^* \mathcal{N}(\mathbf{y}; \mathbf{c}_j^*, \boldsymbol{\Gamma}_j^*)}$ and conditional means $\mathbf{A}_k^* \mathbf{y} + \mathbf{b}_k^*$. But this may lead us back to the usual issue with distances between summary statistics and also we may have to face the label switching issue, not easily handled within ABC procedures.

Recently, developments regarding the Wasserstein distance have emerged (Delon and Desolneux, 2020; Chen et al., 2019), introducing an optimal transport-based distance between Gaussian mixtures, denoted by MW_2 . The L_2 distance between mixtures is also straightforward to compute. Both distances are recalled in supplementary Section S2. We then derive two procedures respectively referred to as GLLiM-MW2-ABC and GLLiM-L2-ABC, writing sometimes GLLiM-D-ABC to include both cases and other distances D .

The semi-automatic ABC extensions that we propose are summarized in Algorithm 1. Algorithm 1 is presented with two simulated data sets, one for training GLLiM and constructing the surrogate posteriors, and one for the ABC procedure itself, but the same data set could be used. For rejection ABC, the selection also requires to fix a threshold ϵ . It is common practice to set ϵ to a quantile of the computed distances. GLLiM then requires to set K , the number of Gaussians in the mixtures. K can be chosen using model selection criteria (see Deleforge et al., 2015). Its precise value is not critical, all the more so if GLLiM is not used for prediction, directly. See details in Section 6.

Algorithm 1 GLLiM-ABC algorithms – Vector and functional variants

- 1: **Inverse operator learning.** Apply GLLiM on a training set $\mathcal{D}_N = \{(\boldsymbol{\theta}_n, \mathbf{y}_n), n \in [N]\}$ to estimate, for any $\mathbf{z} \in \mathcal{Y}$, the K -Gaussian mixture $p_G(\boldsymbol{\theta} \mid \mathbf{z}; \boldsymbol{\phi}_{K,N}^*)$ in (2) as a first approximation of the true posterior $\pi(\boldsymbol{\theta} \mid \mathbf{z})$, where $\boldsymbol{\phi}_{K,N}^*$ does not depend on \mathbf{z} .
 - 2: **Distances computation.** Consider another set $\mathcal{E}_M = \{(\boldsymbol{\theta}_m, \mathbf{z}_m), m \in [M]\}$. For a given observed \mathbf{y} , do one of the following for $m \in [M]$:
 - Vector summary statistics.** (Section 4.1)
 - GLLiM-E-ABC: Compute statistics $s_1(\mathbf{z}_m) = \mathbb{E}_G[\boldsymbol{\theta} \mid \mathbf{z}_m; \boldsymbol{\phi}_{K,N}^*]$ (3).
 - GLLiM-EV-ABC: Compute both $s_1(\mathbf{z}_m)$ and $s_2(\mathbf{z}_m)$ by considering also posterior log-variances, *i.e.* the logarithms of the diagonal elements of (4).
 - In both cases, compute standard distances between summary statistics.
 - Functional summary statistics.** (Section 4.2)
 - GLLiM-MW2-ABC: Compute $MW_2(p_G(\cdot \mid \mathbf{z}_m; \boldsymbol{\phi}_{K,N}^*), p_G(\cdot \mid \mathbf{y}; \boldsymbol{\phi}_{K,N}^*))$.
 - GLLiM-L2-ABC: Compute $L_2(p_G(\cdot \mid \mathbf{z}_m; \boldsymbol{\phi}_{K,N}^*), p_G(\cdot \mid \mathbf{y}; \boldsymbol{\phi}_{K,N}^*))$.
 - 3: **Sample selection.** Select $\boldsymbol{\theta}_m$ values that lead to distances under an ϵ threshold (rejection ABC) or apply an ABC procedure that can handle distances, directly.
 - 4: **Sample use.** For a given observed \mathbf{y} , use the produced sample of $\boldsymbol{\theta}$ values to compute a closer approximation of $\pi(\boldsymbol{\theta} \mid \mathbf{y})$.
-

5 Theoretical properties

Before illustrating GLLiM-D-ABC performance, we investigate the theoretical properties of our ABC quasi-posterior defined via surrogate posteriors.

Let $\mathcal{X} = \Theta \times \mathcal{Y}$ and $(\mathcal{X}, \mathcal{F})$ be a measurable space. Let λ be a σ -finite measure on \mathcal{F} . Whenever we mention below that a probability measure \Pr on \mathcal{F} has a density, we will understand that it has a Radon–Nikodym derivative with respect to λ (λ can typically be chosen as the Lebesgue measure on the Euclidean space). For all $p \in [1, \infty)$ and f, g in appropriate spaces, let $D_p(f, g) = (\int |f(\mathbf{x}) - g(\mathbf{x})|^p d\lambda(\mathbf{x}))^{1/p}$ denote the L_p distance and $D_H^2(f, g) = \int (\sqrt{f(\mathbf{x})} - \sqrt{g(\mathbf{x})})^2 d\lambda(\mathbf{x})$ be the

squared Hellinger distance. When not specified otherwise, let D be an arbitrary distance on \mathcal{Y} or on densities, depending on the context. We further denote the L_p norm for vectors by $\|\cdot\|_p$.

In a GLLiM-D-ABC procedure, the ABC quasi-posterior is constructed as follows. Let $p_G^{K,N}(\boldsymbol{\theta} | \mathbf{y}) = p_G(\boldsymbol{\theta} | \mathbf{y}; \boldsymbol{\phi}_{K,N}^*)$ be the surrogate conditional distribution of form (2), learned from a preliminary GLLiM model with K components and using a learning set $\mathcal{D}_N = \{(\boldsymbol{\theta}_n, \mathbf{y}_n), n \in [N]\}$. This conditional distribution is a K -component mixture, which depends on a set of learned parameters $\boldsymbol{\phi}_{K,N}^*$, independent of \mathbf{y} . The GLLiM-D-ABC quasi-posterior resulting from the GLLiM-D-ABC procedure then depends both on K, N and the tolerance level ϵ and can be written as

$$q_{G,\epsilon}^{K,N}(\boldsymbol{\theta} | \mathbf{y}) \propto \pi(\boldsymbol{\theta}) \int_{\mathcal{Y}} \mathbf{1}_{\{D(p_G^{K,N}(\cdot | \mathbf{y}), p_G^{K,N}(\cdot | \mathbf{z})) \leq \epsilon\}} f_{\boldsymbol{\theta}}(\mathbf{z}) d\mathbf{z}, \quad (5)$$

where D is a distance on densities such as the MW_2 and L_2 metrics, which are both proper distances (see supplementary Section S2).

We provide two types of results, below. In the first result (Theorem 1), the true posterior is used to compare samples \mathbf{y} and \mathbf{z} . This result aims at providing insights on the proposed quasi-posterior formulation and at illustrating its potential advantages. In the second result (Theorem 2), a surrogate posterior is learned and used to compare samples. Conditions are specified under which the resulting ABC quasi-posterior converges to the true posterior.

5.1 Convergence of the ABC quasi-posterior

In this section, we assume a fixed given observed \mathbf{y} and the dependence on \mathbf{y} is omitted from the notation, when there is no confusion.

Let us first recall the standard form of the ABC quasi-posterior, omitting summary statistics from the notation:

$$\pi_{\epsilon}(\boldsymbol{\theta} | \mathbf{y}) \propto \pi(\boldsymbol{\theta}) \int_{\mathcal{Y}} \mathbf{1}_{\{D(\mathbf{y}, \mathbf{z}) \leq \epsilon\}} f_{\boldsymbol{\theta}}(\mathbf{z}) d\mathbf{z}. \quad (6)$$

If D is a distance and $D(\mathbf{y}, \mathbf{z})$ is continuous in \mathbf{z} , the ABC posterior in (6) can be shown to have the desirable property of converging to the true posterior when ϵ tends to 0 (see Prangle et al., 2018).

The proof is based on the fact that when ϵ tends to 0, due to the property of the distance D , the set $\{\mathbf{z} \in \mathcal{Y} : D(\mathbf{y}, \mathbf{z}) \leq \epsilon\}$ in (6) tends to the singleton $\{\mathbf{y}\}$ so that consequently \mathbf{z} in the likelihood can be replaced by the observed \mathbf{y} , which leads to an ABC quasi-posterior proportional to $\pi(\boldsymbol{\theta})f_{\boldsymbol{\theta}}(\mathbf{y})$ and therefore to the true posterior as desired (see also Rubio and Johansen, 2013; Bernton et al., 2019). It is interesting to note that this proof is based on working on the term under the integral only and is using the equality, at convergence, of \mathbf{z} to \mathbf{y} , which is actually a stronger than necessary assumption for the result to hold. Alternatively, if we first rewrite (6) using Bayes' theorem, it follows that

$$\pi_{\epsilon}(\boldsymbol{\theta} | \mathbf{y}) \propto \int_{\mathcal{Y}} \mathbf{1}_{\{D(\mathbf{y}, \mathbf{z}) \leq \epsilon\}} \pi(\boldsymbol{\theta}) f_{\boldsymbol{\theta}}(\mathbf{z}) d\mathbf{z} \propto \int_{\mathcal{Y}} \mathbf{1}_{\{D(\mathbf{y}, \mathbf{z}) \leq \epsilon\}} \pi(\boldsymbol{\theta} | \mathbf{z}) \pi(\mathbf{z}) d\mathbf{z}. \quad (7)$$

That is, when accounting for the normalizing constant:

$$\pi_{\epsilon}(\boldsymbol{\theta} | \mathbf{y}) = \frac{\int_{\mathcal{Y}} \mathbf{1}_{\{D(\mathbf{y}, \mathbf{z}) \leq \epsilon\}} \pi(\boldsymbol{\theta} | \mathbf{z}) \pi(\mathbf{z}) d\mathbf{z}}{\int_{\mathcal{Y}} \mathbf{1}_{\{D(\mathbf{y}, \mathbf{z}) \leq \epsilon\}} \pi(\mathbf{z}) d\mathbf{z}}. \quad (8)$$

Using this equivalent formulation, we can then replace $D(\mathbf{y}, \mathbf{z})$ by $D(\pi(\cdot | \mathbf{y}), \pi(\cdot | \mathbf{z}))$, with D now denoting a distance on densities, and obtain the same convergence result when ϵ tends to 0. More specifically, we can show the following general result. Let us define our ABC quasi-posterior as,

$$q_\epsilon(\boldsymbol{\theta} | \mathbf{y}) \propto \pi(\boldsymbol{\theta}) \int_{\mathcal{Y}} \mathbf{1}_{\{D(\pi(\cdot | \mathbf{y}), \pi(\cdot | \mathbf{z})) \leq \epsilon\}} f_{\boldsymbol{\theta}}(\mathbf{z}) d\mathbf{z},$$

which can be written as

$$q_\epsilon(\boldsymbol{\theta} | \mathbf{y}) = \frac{\int_{\mathcal{Y}} \mathbf{1}_{\{D(\pi(\cdot | \mathbf{y}), \pi(\cdot | \mathbf{z})) \leq \epsilon\}} \pi(\boldsymbol{\theta} | \mathbf{z}) \pi(\mathbf{z}) d\mathbf{z}}{\int_{\mathcal{Y}} \mathbf{1}_{\{D(\pi(\cdot | \mathbf{y}), \pi(\cdot | \mathbf{z})) \leq \epsilon\}} \pi(\mathbf{z}) d\mathbf{z}}. \quad (9)$$

The following theorem shows that $q_\epsilon(\cdot | \mathbf{y})$ converges to $\pi(\cdot | \mathbf{y})$ in total variation, for fixed \mathbf{y} . The proof is detailed in supplementary Section S3.1.

Theorem 1. *For every $\epsilon > 0$, let $A_\epsilon = \{\mathbf{z} \in \mathcal{Y} : D(\pi(\cdot | \mathbf{y}), \pi(\cdot | \mathbf{z})) \leq \epsilon\}$. Assume the following:*

- (A1) $\pi(\boldsymbol{\theta} | \cdot)$ is continuous for all $\boldsymbol{\theta} \in \Theta$, and $\sup_{\boldsymbol{\theta} \in \Theta} \pi(\boldsymbol{\theta} | \mathbf{y}) < \infty$;
- (A2) There exists a $\gamma > 0$ such that $\sup_{\boldsymbol{\theta} \in \Theta} \sup_{\mathbf{z} \in A_\gamma} \pi(\boldsymbol{\theta} | \mathbf{z}) < \infty$;
- (A3) $D(\cdot, \cdot) : \Pi \times \Pi \rightarrow \mathbb{R}_+$ is a metric on the functional class $\Pi = \{\pi(\cdot | \mathbf{y}) : \mathbf{y} \in \mathcal{Y}\}$;
- (A4) $D(\pi(\cdot | \mathbf{y}), \pi(\cdot | \mathbf{z}))$ is continuous, with respect to \mathbf{z} .

Under (A1)–(A4), $q_\epsilon(\cdot | \mathbf{y})$ in (9) converges in total variation to $\pi(\cdot | \mathbf{y})$, for fixed \mathbf{y} , as $\epsilon \rightarrow 0$.

It appears that what is important is not to select \mathbf{z} 's that are close (and at the limit equal) to the observed \mathbf{y} but to choose \mathbf{z} 's so that the posterior $\pi(\cdot | \mathbf{z})$ (the term appearing in the integral in (7)) is close (and at the limit equal) to $\pi(\cdot | \mathbf{y})$. And this last property is less demanding than $\mathbf{z} = \mathbf{y}$. Potentially, there may be several \mathbf{z} 's satisfying $\pi(\cdot | \mathbf{z}) = \pi(\cdot | \mathbf{y})$, but this is not problematic when using (7), while it is problematic when following the standard proof as in Bernton et al. (2019).

5.2 Convergence of the ABC quasi-posterior with surrogate posteriors

In most ABC settings, based on data discrepancy or summary statistics, the above consideration and result are not useful because the true posterior is unknown by construction and cannot be used to compare samples. However this principle becomes useful in our setting, which is based on surrogate posteriors. While the previous result can be seen as an oracle of sorts, it is more interesting in practice to investigate whether a similar result holds when using surrogate posteriors in the ABC likelihood. This is the goal of Theorem 2 below, which we prove for a restricted class of target distribution and of surrogate posteriors that are learned as mixtures.

We now assume that $\mathcal{X} = \Theta \times \mathcal{Y}$ is a compact set and consider the following class $\mathcal{H}_{\mathcal{X}}$ of distributions on \mathcal{X} , $\mathcal{H}_{\mathcal{X}} = \{g_\varphi : \varphi \in \Psi\}$, with constraints on the parameters, Ψ being a bounded parameter set. In addition the densities in $\mathcal{H}_{\mathcal{X}}$ are assumed to satisfy for any $\varphi, \varphi' \in \Psi$, there exists arbitrary positive scalars a, b and B such that

$$\text{for all } \mathbf{x} \in \mathcal{X}, a \leq g_\varphi(\mathbf{x}) \leq b \text{ and } \sup_{\mathbf{x} \in \mathcal{X}} |\log g_\varphi(\mathbf{x}) - \log g_{\varphi'}(\mathbf{x})| \leq B \|\varphi - \varphi'\|_1.$$

We denote by p^K a K -component mixture of distributions from $\mathcal{H}_{\mathcal{X}}$ and defined for all $\mathbf{y} \in \mathcal{Y}$, $p^{K,N}(\cdot | \mathbf{y})$ as follows:

$$\forall \boldsymbol{\theta} \in \Theta, \quad p^{K,N}(\boldsymbol{\theta} | \mathbf{y}) = p^K(\boldsymbol{\theta} | \mathbf{y}; \boldsymbol{\phi}_{K,N}^*),$$

with $\boldsymbol{\phi}_{K,N}^*$ the maximum likelihood estimate (MLE) for the data set $\mathcal{D}_N = \{(\boldsymbol{\theta}_n, \mathbf{y}_n), n \in [N]\}$, generated from the true joint distribution $\pi(\cdot, \cdot)$:

$$\boldsymbol{\phi}_{K,N}^* = \arg \max_{\boldsymbol{\phi} \in \Phi} \sum_{n=1}^N \log(p^K(\boldsymbol{\theta}_n, \mathbf{y}_n; \boldsymbol{\phi})).$$

For every $\epsilon > 0$, let $A_{\epsilon, \mathbf{y}}^{K,N} = \{\mathbf{z} \in \mathcal{Y} : D(p^{K,N}(\cdot | \mathbf{y}), p^{K,N}(\cdot | \mathbf{z})) \leq \epsilon\}$ and $q_{\epsilon}^{K,N}$ denote the ABC quasi-posterior defined with $p^{K,N}$ by

$$q_{\epsilon}^{K,N}(\boldsymbol{\theta} | \mathbf{y}) \propto \pi(\boldsymbol{\theta}) \int_{\mathcal{Y}} \mathbf{1}_{\{D(p^{K,N}(\cdot | \mathbf{y}), p^{K,N}(\cdot | \mathbf{z})) \leq \epsilon\}} f_{\boldsymbol{\theta}}(\mathbf{z}) d\mathbf{z}. \quad (10)$$

Theorem 2. *Assume the following: $\mathcal{X} = \Theta \times \mathcal{Y}$ is a compact set and*

- (B1) *For joint density π , there exists G_{π} a probability measure on Ψ such that, with $g_{\varphi} \in \mathcal{H}_{\mathcal{X}}$, $\pi(\mathbf{x}) = \int_{\Psi} g_{\varphi}(\mathbf{x}) G_{\pi}(d\varphi)$;*
- (B2) *The true posterior density $\pi(\cdot | \cdot)$ is continuous with respect to $\boldsymbol{\theta}$ and \mathbf{y} ;*
- (B3) *$D(\cdot, \cdot) : \Pi \times \Pi \rightarrow \mathbb{R}_+ \cup \{0\}$ is a metric on a functional class Π , which contains the class $\{p^{K,N}(\cdot | \mathbf{y}) : \mathbf{y} \in \mathcal{Y}, K \in \mathbb{N}^*, N \in \mathbb{N}^*\}$. In particular, $D(p^{K,N}(\cdot | \mathbf{y}), p^{K,N}(\cdot | \mathbf{z})) = 0$, if and only if $p^{K,N}(\cdot | \mathbf{y}) = p^{K,N}(\cdot | \mathbf{z})$;*
- (B4) *For every $\mathbf{y} \in \mathcal{Y}$, $\mathbf{z} \mapsto D(p^{K,N}(\cdot | \mathbf{y}), p^{K,N}(\cdot | \mathbf{z}))$ is a continuous function on \mathcal{Y} .*

Then, under (B1)–(B4), the Hellinger distance $D_{\text{H}}(q_{\epsilon}^{K,N}(\cdot | \mathbf{y}), \pi(\cdot | \mathbf{y}))$ converges to 0 in some measure λ , with respect to $\mathbf{y} \in \mathcal{Y}$ and in probability, with respect to the sample $\{(\boldsymbol{\theta}_n, \mathbf{y}_n), n \in [N]\}$. That is, for any $\alpha > 0, \beta > 0$, it holds that

$$\lim_{\epsilon \rightarrow 0, K \rightarrow \infty, N \rightarrow \infty} \Pr(\lambda(\{\mathbf{y} \in \mathcal{Y} : D_{\text{H}}^2(q_{\epsilon}^{K,N}(\cdot | \mathbf{y}), \pi(\cdot | \mathbf{y})) \geq \beta\}) \leq \alpha) = 1. \quad (11)$$

Sketch of the proof of Theorem 2.

For all $\boldsymbol{\theta} \in \Theta, \mathbf{y} \in \mathcal{Y}$, the quasi-posterior (10) can be written equivalently as

$$q_{\epsilon}^{K,N}(\boldsymbol{\theta} | \mathbf{y}) = \int_{\mathcal{Y}} K_{\epsilon}^{K,N}(\mathbf{z}; \mathbf{y}) \pi(\boldsymbol{\theta} | \mathbf{z}) d\mathbf{z},$$

$$\text{with } K_{\epsilon}^{K,N}(\mathbf{z}; \mathbf{y}) = \frac{\mathbf{1}_{\{D(p^{K,N}(\cdot | \mathbf{y}), p^{K,N}(\cdot | \mathbf{z})) \leq \epsilon\}} \pi(\mathbf{z})}{\int_{\mathcal{Y}} \mathbf{1}_{\{D(p^{K,N}(\cdot | \mathbf{y}), p^{K,N}(\cdot | \tilde{\mathbf{z}})) \leq \epsilon\}} \pi(\tilde{\mathbf{z}}) d\tilde{\mathbf{z}}},$$

where $K_{\epsilon}^{K,N}(\cdot; \mathbf{y})$ is a pdf, with respect to $\mathbf{z} \in \mathcal{Y}$, with compact support $A_{\epsilon, \mathbf{y}}^{K,N} \subset \mathcal{Y}$, by definition of $A_{\epsilon, \mathbf{y}}^{K,N}$ and (B4). Using the relationship between Hellinger and L_1 distances (see details in

supplementary Section S3.2 relations (28) and (29)), it then holds that

$$D_H^2(q_{\epsilon}^{K,N}(\cdot | \mathbf{y}), \pi(\cdot | \mathbf{y})) \leq 2D_H(\pi(\cdot | \mathbf{z}_{\epsilon}^{K,N}), \pi(\cdot | \mathbf{y})), \quad (12)$$

where there exists $\mathbf{z}_{\epsilon, \mathbf{y}}^{K,N} \in B_{\epsilon, \mathbf{y}}^{K,N}$ with

$$B_{\epsilon, \mathbf{y}}^{K,N} = \arg \max_{\mathbf{z} \in A_{\epsilon, \mathbf{y}}^{K,N}} D_1(\pi(\cdot | \mathbf{z}), \pi(\cdot | \mathbf{y})).$$

The next step is to bound the right-hand side of (12) using the triangle inequality with respect to the Hellinger distance D_H . Consider the limit point $\mathbf{z}_{0, \mathbf{y}}^{K,N}$ defined as $\mathbf{z}_{0, \mathbf{y}}^{K,N} = \lim_{\epsilon \rightarrow 0} \mathbf{z}_{\epsilon, \mathbf{y}}^{K,N}$. Since for each $\epsilon > 0$, $\mathbf{z}_{\epsilon, \mathbf{y}}^{K,N} \in A_{\epsilon, \mathbf{y}}^{K,N}$ it holds that $\mathbf{z}_{0, \mathbf{y}}^{K,N} \in A_{0, \mathbf{y}}^{K,N}$, where $A_{0, \mathbf{y}}^{K,N} = \bigcap_{\epsilon \in \mathbb{Q}_+} A_{\epsilon, \mathbf{y}}^{K,N}$. By continuity of D , $A_{0, \mathbf{y}}^{K,N} = \{\mathbf{z} \in \mathcal{Y} : D(p^{K,N}(\cdot | \mathbf{z}), p^{K,N}(\cdot | \mathbf{y})) = 0\}$ and $A_{0, \mathbf{y}}^{K,N} = \{\mathbf{z} \in \mathcal{Y} : p^{K,N}(\cdot | \mathbf{z}) = p^{K,N}(\cdot | \mathbf{y})\}$, using (B3). The distance on the right-hand side of (12) can then be decomposed in three parts,

$$\begin{aligned} D_H(\pi(\cdot | \mathbf{z}_{\epsilon, \mathbf{y}}^{K,N}), \pi(\cdot | \mathbf{y})) &\leq D_H(\pi(\cdot | \mathbf{z}_{\epsilon, \mathbf{y}}^{K,N}), \pi(\cdot | \mathbf{z}_{0, \mathbf{y}}^{K,N})) \\ &\quad + D_H(\pi(\cdot | \mathbf{z}_{0, \mathbf{y}}^{K,N}), p^{K,N}(\cdot | \mathbf{y})) \\ &\quad + D_H(p^{K,N}(\cdot | \mathbf{y}), \pi(\cdot | \mathbf{y})). \end{aligned} \quad (13)$$

The first term in the right-hand side can be made close to 0 as ϵ goes to 0 independently of K and N . The two other terms are of the same nature, and the definition of $\mathbf{z}_{0, \mathbf{y}}^{K,N}$ yields $p^{K,N}(\cdot | \mathbf{y}) = p^{K,N}(\cdot | \mathbf{z}_{0, \mathbf{y}}^{K,N})$.

Using that $\pi(\cdot | \cdot)$ is a uniformly continuous function in $(\boldsymbol{\theta}, \mathbf{y})$ on a compact set \mathcal{X} and taking the limit $\epsilon \rightarrow 0$, yields $\lim_{\epsilon \rightarrow 0} D_H^2(\pi(\cdot | \mathbf{z}_{\epsilon, \mathbf{y}}^{K,N}), \pi(\cdot | \mathbf{z}_{0, \mathbf{y}}^{K,N})) = 0$ in measure λ , with respect to $\mathbf{y} \in \mathcal{Y}$. Since this result is true whatever the data set \mathcal{D}_N , it also holds in probability with respect to \mathcal{D}_N . That is, given any $\alpha_1 > 0$, $\beta_1 > 0$, there exists $\epsilon(\alpha_1, \beta_1) > 0$ such that for any $0 < \epsilon < \epsilon(\alpha_1, \beta_1)$,

$$\Pr\left(\lambda\left(\left\{\mathbf{y} \in \mathcal{Y} : D_H^2(\pi(\cdot | \mathbf{z}_{\epsilon, \mathbf{y}}^{K,N}), \pi(\cdot | \mathbf{z}_{0, \mathbf{y}}^{K,N})) \geq \beta_1\right\}\right) \geq \alpha_1\right) = 0.$$

Next, we prove that $D_H^2(\pi(\cdot | \mathbf{z}_{0, \mathbf{y}}^{K,N}), p^{K,N}(\cdot | \mathbf{y}))$ (which is equal to $D_H^2(\pi(\cdot | \mathbf{z}_{0, \mathbf{y}}^{K,N}), p^{K,N}(\cdot | \mathbf{z}_{0, \mathbf{y}}^{K,N}))$) and $D_H^2(p^{K,N}(\cdot | \mathbf{y}), \pi(\cdot | \mathbf{y}))$ both converge to 0 in measure λ , with respect to \mathbf{y} and in probability, with respect to \mathcal{D}_N . Such convergence can be obtained via Rakhlin et al. (2005, Corollary 2.2), and Lemma 2 in supplementary Section S3.3.2 which provides the guarantee that we can choose a measurable function $\mathbf{y} \mapsto \mathbf{z}_{0, \mathbf{y}}^{K,N}$. Equation (11) in Theorem 2 follows from the triangle inequality (13). A detailed proof is provided in supplementary Section S3.2.

Remark.

The GLLiM model involving multivariate unconstrained Gaussian distributions does not satisfy the conditions of Theorem 2 so that $p^{K,N}$ cannot be replaced by $p_G^{K,N}$ in the theorem. However as illustrated in Rakhlin et al. (2005), truncated Gaussian distributions with constrained parameters

can meet the restrictions imposed in the theorem. We are not aware of any more general result involving the MLE of Gaussian mixtures. The GLLiM model could as well be replaced by another model satisfying the conditions of the theorem but for practical applications, this model would need to have computational properties such as the tractability of the estimation of its parameters and needs to be efficient in multivariate and potentially high-dimensional settings.

6 Numerical experiments

Our first two examples are commonly used in the ABC literature and are there to illustrate the flexibility of our method, with an *i.i.d.* observation setting ($R = 100$, $d = 2$, $\ell = 5$) in Section 6.1 and a time series model ($R = 1$, $d = 150$, $\ell = 2$) in Section 6.2. In contrast, the other examples aim at departing from the usual benchmark examples in ABC. We therefore consider settings that exhibit posterior distributions with characteristics such as bimodality and heavy tails. We report an experiment derived from a real application in sound source localization, where the posterior distribution has mass on four 1D manifolds (Section 6.3). Other synthetic examples are described in supplementary Section S4.3. All these other examples are run for a single observation in $d = 10$ dimensions. This choice of dimension is relatively low but corresponds to the dimensions met in practice in some targeted real applications. In particular, we are interested in a real remote sensing inverse problem in planetary science, which is illustrated in Section 6.4.

To circumvent the choice of an arbitrary summary statistic, Fearnhead and Prangle (2012) showed that the best summary statistic, in terms of the minimal quadratic loss, was the posterior mean. This posterior mean is not known and needs to be approximated, *e.g.* by linear regression. In this section, the transformations used for the regression part are $(1, y, y^2, y^3, y^4)$ following the procedure suggested in the **abctools** package (Nunes and Prangle, 2015). We refer to this procedure as semi-automatic ABC. This approach using the posterior mean is further developed in Jiang et al. (2017), where a multilayer perceptron deep neural network regression model is employed. The deep neuronal network with multiple hidden layers considered by Jiang et al. (2017) offers stronger representational power to approximate the posterior mean and hence to learn an informative summary statistic, when compared to linear regression models. Improved results were obtained by Jiang et al. (2017), but we did not compare our approach to their method, except by reporting some of their results when relevant. As our main examples are of relatively small dimension d , we also did not draw comparisons with discrepancy-based ABC techniques such as WABC (Bernton et al., 2019) or classification ABC (Gutmann et al., 2018), which are designed for a more data-rich setting. Discrepancy-based results from Nguyen et al. (2020) are reported when available.

The performance of the four proposed GLLiM-ABC schemes summarized in Algorithm 1 is compared to that of semi-automatic ABC. All reported results are obtained with a simple rejection scheme as per instances implemented in the **abc** R package (Csillery et al., 2012). The other schemes available in the **abc** package have been tested but no notable performance differences were observed. In regards to the final sample thresholding (*i.e.*, choice of ϵ), following common practice, all methods retain samples for which the distance to the observation is under a small (*e.g.* 0.1%) quantile of all computed distances.

The **xLLiM** R package, available on the CRAN, is used to learn a GLLiM model with K components from a set \mathcal{D}_N of N simulations from the true model. The GLLiM implementation uses an isotropic constraint except for the first two examples as specified below. The isotropic GLLiM

involves less parameters than the fully-specified GLLiM and we observed that, in the one observation settings, it provided surrogate posteriors of sufficient quality for the ABC selection scheme. The exact meaning of this constraint can be found in Deleforge et al. (2015); Perthame et al. (2017). Another set of simulated couples $(\boldsymbol{\theta}, \mathbf{y})$ of size M is generally used for the ABC rejection scheme unless otherwise specified.

To visualize posterior samples densities, we use a density estimation procedure based on the `ggplot2` R package with a Gaussian kernel.

6.1 Bivariate Beta model

We illustrate here the case where for each possible values of the parameters it is possible to simulate or observe many (R) *i.i.d.* realizations. The observations to be inverted are also made of R *i.i.d.* realizations but assuming a different number is not a problem.

The bivariate Beta model proposed by Crackel and Flegal (2017) and also used by Nguyen et al. (2020); Jiang et al. (2018) is defined with five positive parameters $\theta_1, \dots, \theta_5$ by letting $v_1 = (u_1 + u_3)/(u_5 + u_4)$ and $v_2 = (u_2 + u_4)/(u_5 + u_3)$, where $u_i \sim \text{Gamma}(\theta_i, 1)$, for $i \in [5]$, and setting $z_1 = v_1/(1 + v_1)$ and $z_2 = v_2/(1 + v_2)$. The likelihood for the bivariate random variable $\mathbf{z}^\top = (z_1, z_2)$ is not available in closed form. The observed sample is generated from the model with true values $(\theta_1, \theta_2, \theta_3, \theta_4, \theta_5) = (1, 1, 1, 1, 1)$. The prior on each parameter is taken to be independent and uniform over interval $[0, 5]$.

We fit a GLLiM model with $K = 100$ for *i.i.d.* data (see Section S2.1 in supplementary material) to a set made of $N = 10^5$ 5-dimensional vectors of parameters, each associated to $R = 100$ *i.i.d.* bivariate observations. For ABC procedures, the tolerance threshold ϵ is set to the 0.05% quantile leading to selected samples of size 50, which matches the experiments of Nguyen et al. (2020); Jiang et al. (2018).

The marginal ABC posterior distributions of parameters $\theta_1, \theta_2, \theta_3, \theta_4$ and θ_5 are displayed in Figure S1 of the supplementary material. Results are qualitatively similar to that of Nguyen et al. (2020); Jiang et al. (2018) which use data discrepancies. Our GLLiM-ABC procedures can be seen as direct alternatives to these latter methods. In contrast, to apply semi-automatic ABC requires summary statistics. In absence of candidate summary statistics, it is suggested by Fearnhead and Prangle (2012) to use evenly-spaced quantiles. For comparison, following Jiang et al. (2018) we apply the semi-automatic procedure on 7 quantiles from the first observed dimension and 7 quantiles from the second. Each simulated data set of size $2 \times R$ is then reduced to 14 quantiles.

Although the use of somewhat arbitrary summary statistics is often problematic, we observe that using 14 quantiles in this case provides satisfying results. Visually (see Figures S1 and S2 in the supplementary material), semi-automatic ABC performs well with modes close to the true parameter values except for θ_4 . The GLLiM mixture appears to provide slightly shifted modes that are better located after an ABC step is added, except for GLLiM-L2-ABC. On this example, the L_2 distance seems to perform poorly. Overall the results are qualitatively similar to that in Jiang et al. (2018).

For a more complete comparison, we also apply the other GLLiM-ABC methods with the 14 quantiles summaries. The standard GLLiM implementation is used with $K = 40$ and no constraint. Our GLLiM-ABC procedures apply as easily in this new setting while the discrepancy-based methods described in Bernton et al. (2019); Jiang et al. (2017); Nguyen et al. (2020) are not designed for this situation. Supplementary Figure S2 shows marginal posteriors for the 5 parameters and 5 procedures. The GLLiM-MW2-ABC procedure based on the MW_2 distance is the best while the one based on L_2

Table 1 Bivariate Beta model: Empirical parameter means, and RMSE for ABC posterior samples averaged over 10 repetitions of the experiment with observed data generated with $\theta = (1, 1, 1, 1, 1)$. The ABC posterior values are computed as empirical values over samples of size 50. Average means closest to 1 and best (lowest) average RMSE values are in boldface. The best results obtained in Nguyen et al. (2020) using various data discrepancies are also given for comparison.

Procedure	$\bar{\theta}_1$	$\bar{\theta}_2$	$\bar{\theta}_3$	$\bar{\theta}_4$	$\bar{\theta}_5$	$R(\theta_1)$	$R(\theta_2)$	$R(\theta_3)$	$R(\theta_4)$	$R(\theta_5)$
GLLiM mixture	2.510	2.546	2.714	2.630	2.591	2.145	2.291	2.201	2.277	2.056
GLLiM-E-ABC	1.439	1.051	0.914	1.095	1.264	0.952	0.791	0.483	0.629	0.510
GLLiM-EV-ABC	1.444	1.037	0.916	1.153	1.205	1.003	0.751	0.556	0.596	0.521
GLLiM-L2-ABC	1.860	2.301	2.430	2.136	2.620	1.268	1.859	2.008	1.536	1.966
GLLiM-MW2-ABC	1.330	1.000	0.8465	1.056	1.159	0.836	0.781	0.458	0.558	0.448
with 14 quantiles as summaries										
Semi-auto ABC	1.235	1.173	0.948	1.000	1.145	0.7601	0.747	0.597	0.599	0.582
GLLiM mixture	0.922	1.139	1.002	0.917	1.040	1.869	1.802	1.286	1.231	0.993
GLLiM-E-ABC	1.209	1.438	1.146	1.071	1.302	0.699	0.880	0.632	0.597	0.659
GLLiM-EV-ABC	1.215	1.565	1.157	1.084	1.167	0.748	0.999	0.677	0.660	0.599
GLLiM-L2-ABC	3.339	2.989	3.420	3.315	2.601	2.711	2.462	2.655	2.715	1.958
GLLiM-MW2-ABC	1.159	1.460	1.146	1.079	1.264	0.687	0.877	0.607	0.593	0.634
Best results using data discrepancies as reported in Nguyen et al. (2020)										
	1.286	1.235	1.083	1.128	1.258	0.828	0.745	0.496	0.498	0.446

performs very poorly. GLLiM-MW2-ABC and GLLiM-E-ABC perform similarly while the addition of log-variances in GLLiM-EV-ABC does not seem to improve posterior estimation.

For a more quantitative comparison, we compute for each posterior samples of size S , empirical means of the parameters, $\bar{\theta}_j = \frac{1}{S} \sum_{i=1}^S \theta_j^i$, and empirical root mean square errors (RMSE) defined as $R(\theta_j) = \sqrt{\frac{1}{S} \sum_{i=1}^S (\theta_j^i - \theta_j^0)^2}$ where $j \in [5]$, $S = 50$ and θ_j^i is the sample i for θ_j and θ_j^0 is the true parameter value. Table 1 shows these quantities averaged over 10 repetitions of the same experiment. The RMSE reported in Table 1 confirm the good performance of semi-automatic ABC when using quantiles as summary statistics and of the GLLiM-MW2-ABC method in both cases, with or without summary statistics. Overall, all methods have similar performance except for GLLiM-L2-ABC. The results also show that there is most of the time a significant gain in running an ABC step after the GLLiM mixture approximation. Since our setting is the same as in Nguyen et al. (2020), we also show in Table 1 the best results obtained for this example as reported in Nguyen et al. (2020). Our results are very similar despite the fact that we use much less realisations with only $R = 100$ *i.i.d.* observations while 500 are used in Nguyen et al. (2020).

6.2 Moving average model

The moving average model is widely used in time series analysis. In particular the moving average model of order 2, MA(2), has often illustrated ABC procedures (Marin et al., 2012; Jiang et al., 2018, 2017; Fearnhead and Prangle, 2012; Nguyen et al., 2020). Natural summary statistics are the empirical auto-covariances of lag 1 and 2. This example is a way to illustrate our method on time series in the same manner as Bernton et al. (2019). To favor comparison with other results on the MA(2) model, we adopt a similar setting as in most papers, *i.e.* that of Jiang et al. (2017), but a quantitative comparison is not strictly possible as the simulated observations may vary from one paper to another. The MA(2) process is a stochastic process $(y'_t)_{t \in \mathbb{N}^*}$ defined by

$$y'_t = z_t + \theta_1 z_{t-1} + \theta_2 z_{t-2}, \quad (14)$$

where $\{z_t\}$ is an *i.i.d.* sequence, according to a standard normal distribution and θ_1 and θ_2 are scalar parameters. A standard identifiability condition is imposed on this model leading to a prior distribution on the triangle described by the inequalities $-2 < \theta_1 < 2$, $\theta_1 + \theta_2 > -1$, $\theta_1 - \theta_2 < 1$. The prior on the two model parameters is taken uniform over the triangular domain. For each pair of parameters (θ_1, θ_2) in the triangular domain, a series of length 150 is simulated according to model (14). This is repeated $N = 10^5$ times. The series to be inverted is simulated similarly with true parameters $\theta_1 = 0.6$ and $\theta_2 = 0.2$. For ABC procedures, the tolerance threshold ϵ is set to the 0.1% quantile leading to selected samples of size 100.

To learn a GLLiM model with $d = 150$, $\ell = 2$, $K = 30$, we propose to use the *i.i.d.* adaptation of GLLiM (see supplementary material S1.2). This corresponds to the approximation suggested in Section 4.2 of Bernton et al. (2019). In terms of GLLiM, this is equivalent to assume bloc diagonal covariance matrices when approximating the likelihood. There is some flexibility as regards the bloc sizes. Larger blocs depart less from the true MA(2) model while requiring more parameters to be estimated. Smaller blocs correspond to neglect some of the dependencies between the blocs but may be acceptable if the remaining dependencies carry enough information on the parameters. Two bloc decompositions are tested. All series of length 150 (y_1, \dots, y_{150}) are first cut into $R = 50$ smaller series of length 3, (y_1, y_2, y_3) , $(y_4, y_5, y_6), \dots$, which are considered as independent and identically distributed. GLLiM is applied with $d = 3$, $R = 50$ and no constraint on the 3×3 blocs themselves. A second experiment is made with $R = 5$ and $d = 30$ *i.e.* with 5 unconstrained blocs of size 30×30 . A better precision especially on θ_2 is obtained with this later setting. This confirms the sensitivity of the dependence over time information in the MA(2) model. We thus choose this setting considering each time series as a sample of 5 smaller series of length 30.

We compare with semi-automatic ABC applied directly to the time series of length 150. Rejection ABC is then also applied using the two empirical auto-covariances as summary statistics. Empirical values for parameter means, standard deviations and correlation, when applying the different ABC schemes for one observed time series, are compared to the true ones computed numerically with importance sampling. The corresponding ABC estimations and samples are shown in supplementary material Table S2 and Figure S3. The results are qualitatively similar to that of Jiang et al. (2017) with a poor estimation of the means for semi-automatic ABC on the full time series. They also confirm results already observed in previous works, namely that semi-automatic and auto-covariance-based procedures do not well capture correlation information between θ_1 and θ_2 .

We then repeat the comparison for 100 different observed series, all simulated from true parameters (0.6, 0.2). In each case, the true posterior means, standard deviations of θ_1 and θ_2 , and correlation are computed numerically. The mean squared errors (MSE) to the true posterior values are then computed and reported in Table 2. The first line in Table 2 shows the averages over the 100 experiments of the posterior true quantities, numerically computed. In particular, we see that the averaged posterior means get close to the true values 0.6 and 0.2. The best MSE are obtained with GLLiM-MW2-ABC while semi-automatic ABC applied directly on the time series provides the largest errors. Semi-automatic ABC provides much lower errors when applied on auto-covariances. The methods using auto-covariances provide satisfying results for the θ_1 mean but not for the other quantities. The GLLiM mixture provides better estimates than semi-automatic ABC on the full time series but remains far from the best performance. This illustrates again that there is a clear gain in complementing GLLiM with an ABC step and that the initial GLLiM mixture needs not to be very accurate. The second best method is GLLiM-E-ABC, which performs similarly as GLLiM-L2-ABC, while surprisingly adding the log-variances in GLLiM-EV-ABC seems to degrade the performance.

Table 2 MA(2) model: mean squared errors (MSE) over 100 simulated observations with the same true parameters (0.6,0.2). MSE are computed for all methods, for the estimated parameter means, standard deviations and correlations compared to their true counterparts computed numerically. The last line shows values as reported in Jiang et al. (2017) based on a deep neural network learning (DNN). The "Exact" line reports the means of the 100 true posterior values. Best (lowest) MSE values are in boldface.

Procedure	mean(θ_1)	mean(θ_2)	std(θ_1)	std(θ_2)	cor(θ_1, θ_2)
	Average				
Exact	0.5986	0.2005	0.0803	0.0815	0.4737
	MSE				
Auto-cov Rejection ABC	0.0055	0.0162	0.0011	0.0080	0.1098
Auto-cov Semi-auto	0.0055	0.0163	0.0013	0.0080	0.1124
Semi-auto ABC	0.3545	0.0239	0.1690	0.1301	0.2592
GLLiM mixture	0.0132	0.0041	0.1649	0.0383	0.1789
GLLiM-E-ABC	0.0033	0.0029	0.0004	0.0004	0.0326
GLLiM-EV-ABC	0.0058	0.0040	0.0027	0.0012	0.0385
GLLiM-L2-ABC	0.0031	0.0034	0.0005	0.0006	0.0340
GLLiM-MW2-ABC	0.0023	0.0018	0.0002	0.0003	0.0256
ABC-DNN Jiang et al. (2017)	0.0096	0.0089	0.0025	0.0026	0.0517

This illustrates the fact that in this unimodal posterior case, the posterior expectation is a good summary statistic. Note however, that GLLiM-MW2-ABC still provides a performance gain. Results with the L_2 distance are much better than in the bivariate Beta example, but not as good as with the MW_2 distance. This may come from the choice of ϵ . We observed in other simulations (supplementary Section S4.5.1) that the L_2 distance was more sensitive to this choice than the MW_2 distance, requiring lower ϵ but this was not further investigated here. To compare with another method that uses estimates of posterior expectations as summary statistics, we report results given in Jiang et al. (2017). Their deep neural network-based method (DNN) provides larger MSE than any of our GLLiM-ABC methods.

6.3 Sound source localization

Our main targets are posterior distributions with multiple modes for which our method is more likely to provide significantly better performance than existing approaches. It is straightforward to construct models that lead to multimodal posteriors by considering likelihoods that are invariant by some transformation. Such non-identifiable models include ill-posed inverse problems that can be constructed as explained in Section S4.3 of the supplementary material. Three synthetic examples therein show that the expectation as a summary statistic suffers from the presence of two equivalent modes, while GLLiM-D-ABC procedures well capture multimodality.

In this sub-section, we consider more complex non-identifiable examples constructed from a real sound source localization problem in audio processing. Although microphone arrays provide the most accurate sound source localization, setups limited to two microphones, *e.g.* Beal et al. (2003); Hospedales and Vijayakumar (2008), are often considered to mimic binaural hearing that resembles the human head with applications such as autonomous humanoid robot modelling.

6.3.1 Two microphone setup

We first consider an artificial two microphone setup in a 2D scene. The object of interest is a sound source located at an unknown position $\theta = (x, y)$. The two microphones are assumed to be located at known positions, respectively denoted by \mathbf{m}_1 and \mathbf{m}_2 . A good cue for the sound source localization

is the interaural time difference (ITD) (Wang and Brown, 2006). The ITD is the difference between two times: the time a sound emitted from the source is acquired by microphone 1 at \mathbf{m}_1 and the time at microphone 2 at \mathbf{m}_2 . The function F that maps a location $\boldsymbol{\theta}$ onto an ITD observation is

$$F(\boldsymbol{\theta}) = \frac{1}{c}(\|\boldsymbol{\theta} - \mathbf{m}_1\|_2 - \|\boldsymbol{\theta} - \mathbf{m}_2\|_2), \quad (15)$$

where c is the sound speed in real applications but set to 1 in our example for the purpose of illustration. The important point is that an ITD value does not correspond to a unique point in the scene space, but rather to a whole surface of points. In fact, each isosurface defined by (15) is represented by one sheet of a two-sheet hyperboloid in 2D. Hence, each ITD observation constrains the location of the auditory source to lie on a 1D manifold. The corresponding hyperboloid is determined by the sign of the ITD. In our example, to create a multimodal posterior, we modify the usual setting by taking the absolute value of the ITD so that solutions can now lie on either of the two hyperboloids. In addition we assume that ITDs are observed with some Student t noise that implies heavy tails and possible outliers. Although the ITD is a univariate measure, we consider a more general d dimensional setting by defining the following Student t likelihood, $\mathbf{y} = (y_1, \dots, y_d)$ and $\text{ITD}(\boldsymbol{\theta}) = |\|\boldsymbol{\theta} - \mathbf{m}_1\|_2 - \|\boldsymbol{\theta} - \mathbf{m}_2\|_2|$, where

$$f_{\boldsymbol{\theta}}(\mathbf{y}) = \mathcal{S}_d(\mathbf{y}; \text{ITD}(\boldsymbol{\theta})\mathbf{I}_d, \sigma^2\mathbf{I}_d, \nu). \quad (16)$$

The above likelihood corresponds to a d -variate Student t -distribution with a d -dimensional location parameter with all dimensions equal to $\text{ITD}(\boldsymbol{\theta})$, diagonal isotropic scale matrix equal to $\sigma^2\mathbf{I}_d$ and degree-of-freedom (dof) parameter ν .

The parameter space is assumed to be $\Theta = [-2, 2] \times [-2, 2]$ and the prior on $\boldsymbol{\theta}$ is assumed to be uniform on Θ . The microphones' positions are $\mathbf{m}_1 = (-0.5, 0)$ and $\mathbf{m}_2 = (0.5, 0)$. We assume $\nu = 3$ and $\sigma^2 = 0.01$. The true $\boldsymbol{\theta}$ is set to $\boldsymbol{\theta} = (1.5, 1)$ and we simulate a 10-dimensional \mathbf{y} following model (16).

The four ABC methods using GLLiM and semi-automatic ABC are compared. Results are reported in supplementary Section S4.4.1.

6.3.2 Two pairs of microphones setting

We build on the previous example to design a more complex setting. Two pairs of microphones are considered respectively located at $((-0.5, 0), (0.5, 0))$ and $((0, -0.5), (0, 0.5))$. The ITD vectors are assumed to be measured with equal probability either from the first pair or from the second pair. It results a likelihood that is a mixture of two equal weight components both following the previous model but for different microphones locations. The 10-dimensional observation \mathbf{y} is generated from a source at location $(1.5, 1)$. Depending on whether this observation is coming from the first pair or second pair component, it results a true posterior as shown in Figure 1 (d) or one with non-intersecting hyperbolas. The contour plot indicates that the observation corresponds to the $((0, -0.5), (0, 0.5))$ pair.

The GLLiM model used consists of $K = 20$ Gaussian components with an isotropic constraint. A selected sample of 1000 values is retained by thresholding the distances under the 0.1% quantile. In a first test, semi-automatic ABC and GLLiM use the same data set of size $M = 10^6$ which is also used for the rejection ABC part. Selected samples are shown in supplementary Section S4.4.2, Figure S8. The mixture provided by GLLiM as an approximation of the true posterior (Figure S8 (d)) well

captures the main posterior parts. This GLLiM posterior is a 20-component Gaussian mixture of form (2). The true posterior expectations are all zero and are thus not informative about the location parameters. However, a correct structure can be seen in the GLLiM-E-ABC sample, in contrast to the semi-automatic one that shows no structure as expected. Adding the posterior log-variance estimations has a good impact on the selected sample, which is only marginally different from the GLLiM-D-ABC samples. This suggests that the posterior log-variances are very informative on the location parameters.

When GLLiM is first learned with a smaller data set of size $N = 10^5$ and different from the rejection ABC data set, results slightly degrade, but not significantly so (Figure 1). More badly localized estimations can be seen in the samples of Figure 1 (a,b), but the GLLiM-D-ABC samples are well localized and are not really impacted by this difference in the GLLiM learning step. In this case the improvement of GLLiM-D-ABC over GLLiM-EV-ABC is clearer.

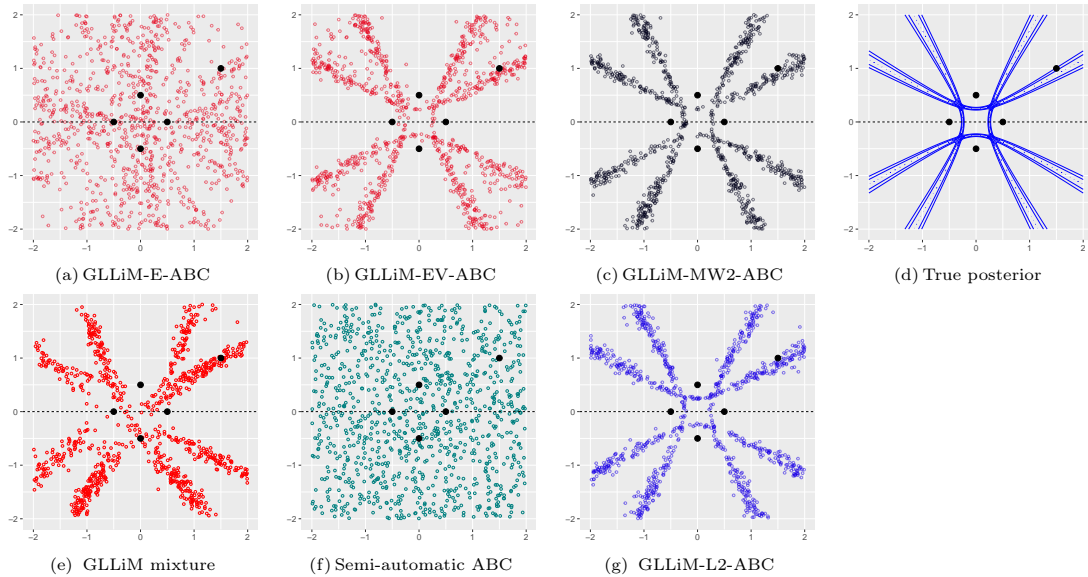


Figure 1 Sound source localization with a mixture of two microphones pairs. GLLiM is learned on a first data set of size $N = 10^5$ while ABC is run using the largest data set of size $M = 10^6$. Selected samples using (a) GLLiM posterior expectations, (b) GLLiM posterior expectations and log variances, (c) MW_2 distances, (d) contours of the true posterior distribution, (e) approximate GLLiM posterior for the observed data, (f) semi-automatic ABC, and (g) L_2 distances. Black points on the dotted line are the microphones positions. The fifth black point is the true sound source localization.

6.4 A physical model inversion in planetary science

As a real-world example, we consider a remote sensing application coming from the study of planetary environment; in particular, the morphological, compositional, photometrical and textural characterization of sites on the surface of a planet. The composition of the surface materials is generally established on the basis of spectral mixing and physical modelling techniques using images produced by hyperspectral cameras, from different angles during a site flyover. An example for the Mars planet

is described by Murchie et al. (2009); Fernando et al. (2016). Such observations can also be measured in the laboratory, on known materials to validate a model. In both cases, the interpretation of the surface Bidirectional Reflectance Distribution Factor (BRDF) extracted from these observations is based on the inversion of a model of radiative transfer, linking physical and observable parameters in a non-linear way.

The Hapke model is a semi-empirical photometric model that relates physically meaningful parameters to the reflectivity of a granular material for a given geometry of illumination and viewing. Formally, it links a set of parameters $\boldsymbol{\theta} \in \mathbb{R}^4$ to a *theoretical* BRDF denoted by $\mathbf{y} = F_{\text{Hapke}}(\boldsymbol{\theta}) \in \mathbb{R}^d$. A given experiment defines d geometries of measurement, each parameterized by a triplet (θ_0, θ, ϕ) of incidence, emergence and azimuth angles. Moreover, $\boldsymbol{\theta} = (\omega, \bar{\theta}, b, c)$ are the sensitive parameters, respectively single scattering albedo, macroscopic roughness, asymmetry parameter and backscattering fraction. More details on these quantities and their photometric meanings may be found in Schmidt and Fernando (2015); Labarre (2017). Although available, the expression of F_{Hapke} is very complex and tedious to handle analytically, with a number of approximations required (see the description of the function in more than 15 pages in Labarre 2017). In practice, it is therefore mainly used via a numerical code, allowing simulations from the model. In addition, previous studies (Kugler et al. 2021; Schmidt and Fernando 2015) have shown evidence for the existence of multiple solutions or for the possibility to obtain very similar observations from different sets of parameters, which makes this setting appropriate for testing the ability of our procedures to recover multimodal posterior distributions.

In the following experiments, all parameters are transformed to be in $[0, 1]^4$, which amounts to keep b and c unchanged, divide $\bar{\theta}$ by 30 and operate the following change of variable for ω , $\gamma = 1 - \sqrt{1 - \omega}$. This last transformation also has the advantage of avoiding the non-linearity of F_{Hapke} , when ω tends to 1. The experimental setting defines geometries at which the measurements are made, which in turn define F_{Hapke} . The number of geometries thus corresponds to the size d , of each observation. The measurement geometries used to define F_{Hapke} are borrowed from a real laboratory experiment presented below. The number of parameters is therefore $\ell = 4$ with $d = 10$ observed geometries. The sets to learn GLLiM and generate ABC samples are both set to size $N = M = 10^5$. For each couple $(\boldsymbol{\theta}, \mathbf{y})$ in the simulated data sets, the 4 parameters $(\boldsymbol{\theta})$ are simulated uniformly in $[0, 1]^4$. Following a previous study (Kugler et al., 2021), the corresponding reflectance curves are generated as $\mathbf{y} = F_{\text{Hapke}}(\boldsymbol{\theta}) + \boldsymbol{\varepsilon}$, where $\boldsymbol{\varepsilon}$ is a centered Gaussian variable with isotropic covariance $\sigma^2 \mathbf{I}_d$. In this section $\sigma = 0.05$. The GLLiM model is learned with $K = 40$.

Prior to real data inversion, performance is assessed by considering an observation simulated from the Hapke model, as explained in the supplementary Section S4.5.1. In this experiment, ϵ is varying to observe the behavior of the different methods (Figure S9). GLLiM-L2-ABC seems less robust, than the other procedures, to these variations and even degrades in performance when ϵ is too high. The two procedures based on expectations show satisfying performance with globally less sharp posteriors. The addition of the posterior log-variances does not seem to significantly change the selected samples.

Reflectance measurements made in the laboratory are also generally considered by experts (see *e.g.* Pilorget et al. 2016). We focus on one observation coming from a mineral called Nontronite (see Kugler et al. 2021 for a description). The experiment consists of taking measures at 100 wavelengths in the spectral range 400–2800 nm. Each of these 100 measures is an observation to be inverted. We focus on one of them, at 2310 nm. This observation has been chosen from previous study (Kugler et al., 2021) as likely to exhibit multiple solutions. The size d of each observation is $d = 10$ and

the corresponding angles are such that the incidence and azimuth angles are fixed to $\theta_0 = 45$ and $\phi = 0$. This number d of geometries is typical of real observations for which the number of possible measurements during a planet flyover is limited.

Figure 2 provides the posterior marginals for the Nontronite, obtained by setting ϵ to the 0.1% quantile of the distances. Two solutions can be deduced. Parameters ω and c show unimodal posterior distributions while $\bar{\theta}$ distribution exhibits two modes. For b , the GLLiM-MW2-ABC sample shows a second smaller mode around 0.5 but this mode is not maintained when ϵ is set to a lower quantile (see Figure S10 in supplementary Section S4.5.2). We therefore consider that the multiplicity comes mainly from $\bar{\theta}$. In the absence of ground truth, it is difficult to fully validate the estimations. However a simple inspection consists of checking the reconstructed signals. The top-right plot in Figure 2 compares the inverted signal to the reconstructed signals obtained by applying the Hapke model to the two sets of estimated parameters, namely $(0.59, 0.15, 0.14, 0.06)$ and $(0.59, 0.42, 0.14, 0.06)$, which differ only in $\bar{\theta}$. The proximity of the reconstructions confirms the existence of multiple solutions and thus the relevance of a multimodal posterior. One solution can be selected by choosing the parameters that provides the best reconstruction. The set $(0.59, 0.42, 0.14, 0.06)$ is selected as its MSE is slightly lower (2.6×10^{-4} vs 3.3×10^{-4}). This is satisfactory, as the lower value of $\bar{\theta}$ in the other solution is less physically interpretable. Note that for simplicity, we have used a uniform prior on θ but for a more meaningful study in planetary science, information on the parameters plausible values could be incorporated directly in the prior.

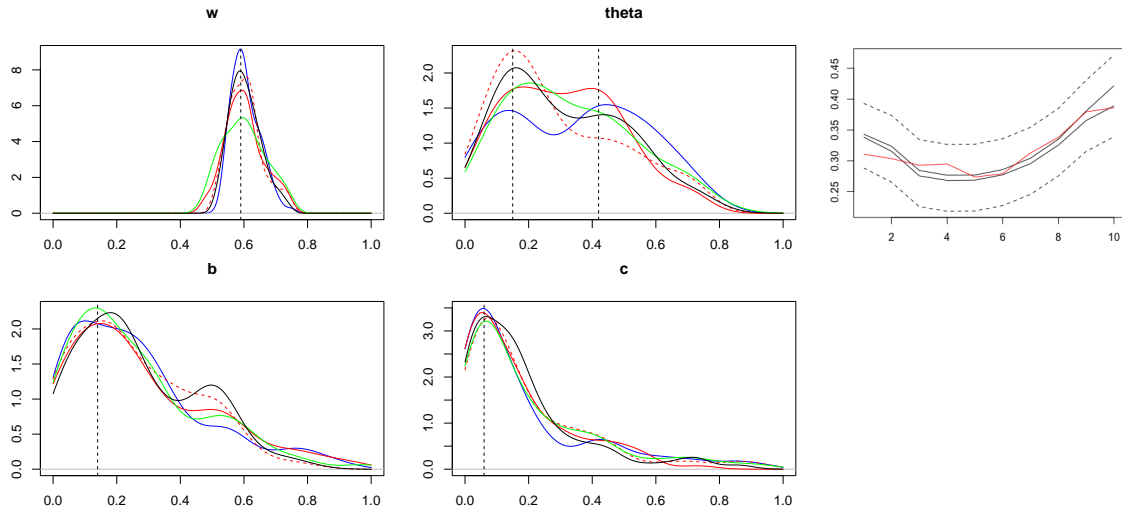


Figure 2 Real observation inversion using the Hapke model. Posterior marginals for $\omega, \bar{\theta}, b$ and c with GLLiM-E-ABC (red), GLLiM-EV-ABC (dotted red), semi-automatic ABC (green), GLLiM-L2-ABC (blue) and GLLiM-MW2-ABC (black). The threshold ϵ is set to the 0.1% quantile (100 selected values). The vertical lines indicate the values $(\omega, \bar{\theta}, b, c) = (0.59, 0.15, 0.14, 0.06)$ and $(0.59, 0.42, 0.14, 0.06)$. The corresponding signal reconstructions (black lines) are shown in the top-right plot with the observed signal in red. The dashed lines correspond to the addition/subtraction of a standard deviation of 0.05 around the reconstructions.

7 Conclusion and perspectives

In this work, the issue of choosing summary statistics was revisited. We built on the seminal work of Fearnhead and Prangle (2012) and their semi-automatic ABC by replacing the approximate posterior expectations with functional statistics; namely approximations of the posterior distributions. These surrogate posterior distributions were obtained in a preliminary learning step, based on an inverse regression principle. This is original with respect to most standard regression procedures, which usually provide only point-wise predictions, *i.e.* first order moments. So doing, we not only could compute approximate posterior moments of higher orders as summary statistics but, more generally, approximate full posterior distributions. This learning step was based on the so-called GLLiM model, which provides surrogate posteriors in the parametric family of Gaussian mixtures. Preliminary experiments showed that although the posterior moments provided by GLLiM were not always leading to better results than that provided by semi-automatic ABC, the use of the full surrogate posteriors was always an improvement. Consequently, an interesting feature of our approach is that, with our adaptation of the original GLLiM model to *i.i.d.* data, it can be seen as an alternative to both summary-based and discrepancy-based procedures.

To handle distributions as functional summary statistics, our procedure required appropriate distances. We investigated an L_2 and a Wassertein-based distance (MW_2). The two distances often performed similarly but poor results have been observed with L_2 that would require further investigations. The MW_2 distance appeared to be more robust. As illustrated in our remote sensing example, it may also allow for the ability to set the tolerance level at a higher value without overly degrading the quality of the posterior sample.

Among aspects that have not been thoroughly investigated in this work, we could refine the way to choose this tolerance level ϵ or combine GLLiM with more sophisticated ABC schemes than the simple rejection scheme.

Another interesting perspective would be to investigate the use of GLLiM in the context of synthetic likelihood (SL) approaches. When used in a Bayesian framework, SL techniques can be viewed as alternatives to ABC in which the intractable likelihood is replaced by an estimator of the likelihood (Price et al., 2018). Since the seminal work of Wood (2010), several estimators have been proposed (e.g. Ong et al., 2018; An et al., 2019, 2020; Frazier and Drovandi, 2021), often derived from auxiliary models (Drovandi et al., 2015). In the ABC framework of this paper, GLLiM was used to provide approximate posteriors but these posteriors are themselves coming from approximate likelihoods that could lead to new SL procedures.

At last, in principle, any other method that is able to provide approximate surrogate posteriors could be used in place of GLLiM to produce the functional summaries. Besides the family of mixture of experts models which are similar to GLLiM, mixture density networks (Bishop, 1994) or normalizing flows (Dinh et al., 2015; Kobyzev et al., 2020; Kruse et al., 2021) are potential candidates. To our knowledge, other common neural networks, like most regression techniques, would not be appropriate as they only focus on point-wise predictions.

Acknowledgements.

FF would like to thank Guillaume Kon Kam King for an initial discussion on semi-automatic ABC, which inspired this work, Benoit Kugler and Sylvain Douté for providing the simulations for the planetary science example and for helpful discussions on the Hapke model.

References

- An, Z., D.J. Nott, and C. Drovandi. 2020. Robust Bayesian Synthetic Likelihood via a Semi-Parametric Approach. *Statistics and Computing* (3): 543–557 .
- An, Z., L. South, D.J. Nott, and C. Drovandi. 2019. Accelerating Bayesian Synthetic Likelihood With the Graphical Lasso. *Journal of Computational and Graphical Statistics* (28): 471–475 .
- Arridge, S., P. Maass, O. Öktem, and C.B. Schönlieb. 2019, May. Solving inverse problems using data-driven models. *Acta Numerica* 28: 1–174 .
- Beal, M.J., N. Jojic, and H. Attias. 2003, July. A graphical model for audiovisual object tracking. *IEEE Trans. Pattern Anal. Mach. Intell.* 25(7): 828–836 .
- Bernard-Michel, C., S. Douté, M. Fauvel, L. Gardes, and S. Girard. 2009. Retrieval of Mars surface physical properties from OMEGA hyperspectral images using Regularized Sliced Inverse Regression. *Journal of Geophysical Research: Planets* 114(E6) .
- Bernton, E., P.E. Jacob, M. Gerber, and C.P. Robert. 2019. Approximate Bayesian computation with the Wasserstein distance. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 81: 235–269 .
- Bishop, C.M. 1994. Mixture density networks. Technical report, Aston University.
- Blum, M.G.B., M.A. Nunes, D. Prangle, and S.A. Sisson. 2013, May. A comparative review of dimension reduction methods in approximate Bayesian computation. *Statistical Science* 28(2): 189–208 .
- Buchholz, A. and N. Chopin. 2019. Improving Approximate Bayesian Computation via Quasi-Monte Carlo. *Journal of Computational and Graphical Statistics* 28(1): 205–219 .
- Chen, Y., T.T. Georgiou, and A. Tannenbaum. 2019. Optimal Transport for Gaussian Mixture Models. *IEEE Access* 7: 6269–6278 .
- Cook, R.D. and L. Forzani. 2019, April. Partial least squares prediction in high-dimensional regression. *The Annals of Statistics* 47(2): 884–908 .
- Crackel, R. and J. Flegal. 2017. Bayesian inference for a flexible class of bivariate beta distributions. *Journal of Statistical Computation and Simulation* 87: 295–312 .
- Csillery, K., O. Francois, and M. Blum. 2012. abc: an R package for approximate Bayesian computation (ABC). *Methods in Ecology and Evolution* .
- Del Moral, P., A. Doucet, and A. Jasra. 2012, September. An Adaptive Sequential Monte Carlo Method for Approximate Bayesian Computation. *Statistics and Computing* 22(5): 1009–1020 .
- Deleforge, A., F. Forbes, S. Ba, and R. Horaud. 2015, September. Hyper-Spectral Image Analysis with Partially-Latent Regression and Spatial Markov Dependencies. *IEEE Journal of Selected*

Topics in Signal Processing 9(6): 1037–1048 .

Deleforge, A., F. Forbes, and R. Horaud. 2015, September. High-Dimensional Regression with Gaussian Mixtures and Partially-Latent Response Variables. *Statistics and Computing* 25(5): 893–911 .

Delon, J. and A. Desolneux. 2020. A Wasserstein-type distance in the space of Gaussian Mixture Models. *SIAM Journal on Imaging Sciences* .

Dinh, L., D. Krueger, and Y. Bengio 2015. NICE: non-linear independent components estimation. In Y. Bengio and Y. LeCun (Eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Workshop Track Proceedings*.

Drovandi, C., T. Pettitt, and A. Lee. 2015. Bayesian indirect inference using a parametric auxiliary model. *Statistical Science* 30(1): 72–95 .

Drovandi, C.C. and A.N. Pettitt. 2011. Likelihood-free Bayesian estimation of multivariate quantile distributions. *Computational Statistics and Data Analysis* 55: 2541–2556 .

Fearnhead, P. and D. Prangle. 2012. 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 .

Fernando, J., F. Schmidt, and S. Douté. 2016. Martian surface microtexture from orbital CRISM multi-angular observations: A new perspective for the characterization of the geological processes. *Planetary and Space Science* 128: 30–51 .

Frazier, D.T. and C. Drovandi. 2021. Robust Approximate Bayesian Inference With Synthetic Likelihood. *Journal of Computational and Graphical Statistics*: 1–19 .

Gutmann, M.U., R. Dutta, S. Kaski, and J. Corander. 2018. Likelihood-free inference via classification. *Statistics and Computing* 28: 411–425 .

Hospedales, T.M. and S. Vijayakumar. 2008, December. Structure inference for Bayesian multi-sensory scene understanding. *IEEE Trans. Pattern Anal. Mach. Intell.* 30(12): 2140–2157 .

Hovorka, R., V. Canonico, L.J. Chassin, U. Haueter, M. Massi-Benedetti, M.O. Federici, T.R. Pieber, H.C. Schaller, L. Schaupp, T. Vering, and M.E. Wilinska. 2004. Nonlinear model predictive control of glucose concentration in subjects with type 1 diabetes. *Physiological Measurement* 25(4): 905–920 .

Ingrassia, S., S.C. Minotti, and G. Vittadini. 2012. Local Statistical Modeling via a Cluster-Weighted Approach with Elliptical Distributions. *Journal of classification* 29(3): 363–401 .

Jiang, B., T.Y. Wu, Z. C., and W. Wong. 2017. Learning summary statistics for Approximate Bayesian Computation via Deep Neural Network. *Statistica Sinica*: 1595–1618 .

- Jiang, B., T.Y. Wu, and W.H. Wong 2018. Approximate Bayesian computation with Kullback-Leibler divergence as data discrepancy. In *21st International Conference on Artificial Intelligence and Statistics (AISTATS)*.
- Kobyzev, I., S. Prince, and M. Brubaker. 2020. Normalizing Flows: An Introduction and Review of Current Methods. *IEEE Trans. Pattern Anal. Mach. Intell.*: 1–1 .
- Kristan, M., A. Leonardis, and D. Skočaj. 2011. Multivariate online kernel density estimation with Gaussian kernels. *Pattern Recognition* 44(10-11): 2630–2642 .
- Kruse, J., L. Ardizzone, C. Rother, and U. Kothe. 2021. Benchmarking invertible architectures on inverse problems. *Workshop on Invertible Neural Networks and Normalizing Flows (ICML 2019)*, *arXiv preprint arXiv:2101.10763* .
- Kugler, B., F. Forbes, and S. Douté. 2021. Fast Bayesian Inversion for high dimensional inverse problems. *To appear in Statistics and Computing*, <https://hal.archives-ouvertes.fr/hal-02908364> .
- Labarre, S. 2017. *Caractérisation et modélisation de la rugosité multi-échelle des surfaces naturelles par télédétection dans le domaine solaire*. Ph. D. thesis, Physique Univers Sorbonne Paris Cité. Supervised by C. Ferrari and S. Jacquemoud.
- Lemasson, B., N. Pannetier, N. Coquery, L.S.B. Boisserand, N. Collomb, N. Schuff, M. Moseley, G. Zaharchuk, E.L. Barbier, and T. Christen. 2016, November. MR Vascular Fingerprinting in Stroke and Brain Tumors Models. *Scientific Reports* 6: 37071 .
- Li, K.C. 1991, June. Sliced Inverse Regression for Dimension Reduction. *Journal of American Statistical Association* 86(414): 316–327 .
- Ma, D., V. Gulani, N. Seiberlich, K. Liu, J.L. Sunshine, J.L. Duerk, and M.A. Griswold. 2013, March. Magnetic Resonance Fingerprinting. *Nature* 495(7440): 187–192 .
- Marin, J.M., P. Pudlo, C.P. Robert, and R.J. Ryder. 2012. Approximate Bayesian computation methods. *Statistics and Computing* 22: 1167–1180 .
- Mesejo, P., S. Sallet, O. David, C. Bénar, J.M. Warnking, and F. Forbes. 2016, March. A differential evolution-based approach for fitting a nonlinear biophysical model to fMRI BOLD data. *IEEE Journal of Selected Topics in Signal Processing* 10(2): 416–427 .
- Muandet, K., K. Fukumizu, F. Dinuzzo, and B. Scholkopf 2012. Learning from distributions via support measure machines. In *Advances in neural information processing systems*, pp. 10–18.
- Murchie, S.L., F.P. Seelos, C.D. Hash, D.C. Humm, E. Malaret, J.A. McGovern, T.H. Choo, K.D. Seelos, D.L. Buczkowski, M.F. Morgan, O.S. Barnouin-Jha, H. Nair, H.W. Taylor, G.W. Patterson, C.A. Harvel, J.F. Mustard, R.E. Arvidson, P. McGuire, M.D. Smith, M.J. Wolff, T.N. Titus, J.P. Bibring, and F. Poulet. 2009. Compact Reconnaissance Imaging Spectrometer for Mars investigation and data set from the Mars Reconnaissance Orbiter’s primary science phase. *Journal of Geophysical Research: Planets* 114(E2): E00D07 .

- Nataraj, G., J.F. Nielsen, C. Scott, and J.A. Fessler. 2018, September. Dictionary-Free MRI PERK: Parameter Estimation via Regression with Kernels. *IEEE Trans. Med. Imaging* 37(9): 2103–2114 .
- Nguyen, H.D., J. Arbel, H. Lü, and F. Forbes. 2020. Approximate Bayesian Computation Via the Energy Statistic. *IEEE Access* 8: 131683–131698 .
- Nguyen, H.D., F. Chamroukhi, and F. Forbes. 2019, August. Approximation results regarding the multiple-output Gaussian gated mixture of linear experts model. *Neurocomputing* .
- Nguyen, T.T., F. Chamroukhi, H.D. Nguyen, and G.J. McLachlan. 2020. Approximation of probability density functions via location-scale finite mixtures in Lebesgue spaces. *arXiv preprint arXiv:2008.09787* .
- Nunes, M.A. and D. Prangle 2015. abctools: An R package for tuning Approximate Bayesian Computation analyses. <https://cran.r-project.org/web/packages/abctools/>.
- Ong, V., D. Nott, M.N. Tran, S. Sisson, and C. Drovandi. 2018. Likelihood-free inference in high dimensions with synthetic likelihood. *Computational Statistics and Data Analysis* 128 .
- Park, M., W. Jitkrittum, and D. Sejdinovic 2016. K2-ABC: approximate Bayesian computation with kernel embeddings. In *19th International Conference on Artificial Intelligence and Statistics (AISTATS)*.
- Perthame, E., F. Forbes, A. Deleforge, E. Devijver, and M. Gallopin 2017. xLLiM: An R package for High Dimensional Locally-Linear Mapping. <https://cran.r-project.org/web/packages/xLLiM/>.
- Pilorget, C., J. Fernando, B.L. Ehlmann, F. Schmidt, and T. Hiroi. 2016. Wavelength dependence of scattering properties in the VIS–NIR and links with grain-scale physical and compositional properties. *Icarus* 267: 296–314 .
- Prangle, D., R.G. Everitt, and T. Kypraios. 2018. A rare event approach to high-dimensional approximate Bayesian computation. *Statistics and Computing* 28: 819–834 .
- Price, L.F., C.C. Drovandi, A. Lee, and D.J. Nott. 2018. Bayesian Synthetic Likelihood. *Journal of Computational and Graphical Statistics* 27(1): 1–11 .
- Rakhlin, A., D. Panchenko, and S. Mukherjee. 2005. Risk bounds for mixture density estimation. *ESAIM: Probability and Statistics* 9: 220–229 .
- Rubio, F. and A.M. Johansen. 2013. A simple approach to maximum intractable likelihood estimation. *Electronic Journal of Statistics* 7: 1632–1654 .
- Schmidt, F. and J. Fernando. 2015. Realistic uncertainties on Hapke model parameters from photometric measurements. *Icarus* 260: 73–93 .
- Sisson, S.A., Y. Fan, and M.A. Beaumont eds. 2019. *Handbook of Approximate Bayesian Computation*. Boca Raton: CRC Press.

- Soubeyrand, S., F. Carpentier, F. Guiton, and E.K. Klein. 2013. Approximate Bayesian computation with functional statistics. *Statistical Applications in Genetics and Molecular Biology* 12(1): 17–37 .
- Sriperumbudur, B.K., A. Gretton, K. Fukumizu, B. Scholkopf, and G.R. Lanckriet. 2010. Hilbert space embeddings and metrics on probability measures. *The Journal of Machine Learning Research* 11: 1517–1561 .
- Wang, D. and G.J. Brown. 2006. *Computational Auditory Scene Analysis: Principles, Algorithms, and Applications*. Wiley-IEEE Press.
- Wang, F., T. Syeda-Mahmood, B.C. Vemuri, D. Beymer, and A. Rangarajan 2009. Closed-form Jensen-Renyi divergence for mixture of Gaussians and applications to group-wise shape registration. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 648–655. Springer.
- Wiqvist, S., P.A. Mattei, U. Picchini, and J. Frellsen 2019, 09–15 Jun. Partially exchangeable networks and architectures for learning summary statistics in approximate Bayesian computation. In K. Chaudhuri and R. Salakhutdinov (Eds.), *Proceedings of the 36th International Conference on Machine Learning*, Volume 97, Long Beach, California, USA, pp. 6798–6807.
- Wood, S. 2010. Statistical inference for noisy nonlinear ecological dynamic systems. *Nature* 466(7310): 1102–1104 .