

Response to Referee Comments: Parallel Adaptive Importance Sampling

March 5, 2018

First of all, we would like to thank the referees and the associate editor for their time in reading and commenting on our manuscript. In this document we will address each of the points raised, and describe any relevant changes that have been made to the paper. We have made some substantial changes, not least with the inclusion of a new “large scale” example where the likelihood requires a solve of a differential equation. We demonstrate in this example that PAIS is 2 orders of magnitude faster, in terms of CPU time, than the equivalent standard MCMC. We include a copy of the manuscript with all of the changes marked.

Referee 1 comments

Comment: *The PAIS algorithm does exhibit interesting ideas to combine the information obtained from an ensemble of proposals. Furthermore, it seems quite flexible since several proposals and re-sampling methods can be used.*

However, the examples somehow fail to convince that the algorithm can be used in a wide range of challenging inverse problems. The approach is limited by the dimensionality of the problem..

Answer: We thank the referee for their comments. We are very open about the fact that this approach cannot be used for high dimensional problems, but there are a plethora of lower dimensional inverse problems which exhibit structure in the posterior, such as multimodality or high correlations, for which existing MCMC methods struggle to converge in reasonable time. Dimensionality is not the only challenge faced in inverse problems, and these methods have been developed to tackle problems of a different nature. In particular, this method was designed with a view to sampling from inverse problems arising in Biology, where challenging posteriors exist on relatively low dimensional state spaces, and for which the likelihood can be very expensive to compute. We have added some words on this to clarify further the intended type of problem that this method is in the introduction, and again in the discussion at the end.

Comment: *For mixture problems, if the mixture has too many elements (or the number of elements is unknown), then this scheme will have difficulties..*

Answer: This is true, as would state-of-the-art Riemannian manifold MALA, HMC and a whole range of other methods. I would be happy for the referee to indicate a method which would comfortably deal with a mixture problem with a large number of possible known/unknown components. The example was chosen to demonstrate the remarkable performance of the algorithm (with respect to standard MCMC methods) for multimodal posteriors in low dimensions, and we believe it has achieved that. We have added a paragraph to the end of the example to discuss the point that the referee makes.

Comment: *Also, it is not clear how it will behave in real life inverse problems where the likelihood is determined by complex physical systems that involve, for example, ODE's or PDE's. If the likelihood is expensive to evaluate, the tuning of the scaling parameters for the PAIS-RW might be a concern when the number of chains is big..*

Answer: The tuning of the scaling parameters is actually very fast, often converging within a few percentage points of the optimum within the order of 10^2 iterations. We have included some plots and discussion to the end of Section 4.2 to demonstrate this point, and we thank the referee for highlighting this issue. This is quicker than one might expect for an adaptive serial algorithm.

Comment: *An real data example for the PAIS algorithm is much needed. It will be interesting to see how PAIS is compared not only against a naive parallelization, but also compared against a state of the art MCMC algorithm or/and a efficient ad hoc MCMC algorithm designed for such example..*

Answer: Unfortunately we do not have the resources to apply our method to a full-scale real data example. However, following additional comments from the editor, we have included a new example where the likelihood requires the solution of an ODE. In Section 7.5 we include an example of data assimilation for the Lorenz '63 equations, which is a set of equations which is commonly used to test methods for inference in the applied mathematics community.

Expensive likelihoods actually make this method more applicable, since they make the extra overheads of the resampling step relatively less important. This is precisely what the method is designed for, and is the very reason for us to compare the algorithms by looking at error as a function of the number of samples/likelihood evaluations. This is seen even for a simple 3D ODE as we see in our example in Section 7.5, where PAIS reaches a relative error in the mean of 10^{-5} in a time which is 2 orders of magnitude quicker in CPU time than for trivially parallelised RWMH chains.

Comment: *There aren't details about the architecture used to run the M*

chains in parallel. This is relevant since there are several options such as: multi-core and multi-processor computers, GPU cards, clusters, cloud computing, etc. Comments about the advantages and disadvantages of the architecture are also welcome. The same goes for details regarding the programming language used..

Answer: We thank the referee for this comment, as we should have included these details in the manuscript. The method was implemented in C++, serially, and was run on a single core of a Dell server with four 8-core 3.3GHz CPUs, 64Gb memory, and these details have been added in Section 7.

We are not computer scientists, and as such, we are not in a position to comment on specific architectures, especially as the optimal architecture would be problem specific. The optimal architecture would likely depend on the relative costs of communication of the states across the processors, the cost of the likelihood evaluation, and the cost of the resampler for the chosen ensemble size. We have added a discussion to Section 3, following the presentation of the PAIS algorithm.

Comment: *CPU time must be reported to compare the gain of each algorithm in order to have a fair comparison. Authors claim that there is a significant gain in terms of the number of iterations needed to achieve a desire tolerance; examples in sections 7.2 and 7.3 show that PAIS saves up to 90% of the iterations. However, if the CPU time needed to run 10 PAIS iterations is ten times the time needed to run 100 RWMH iterations, then the gain is not as good as presented. PAIS iterations are more time consuming than RWMH iterations due to the re-sampling step. The authors mentioned that the re-sampling step is computationally expensive and that is why they proposed the AMR re-sampling method. CPU time is of particular importance in example 7.4 since the number of parallel chains grows from 50 to 500..*

Answer: The referee may have missed Figure 5.1 in which we show the computational cost, in time, of the resamplers that we use in the paper, both the optimal transport resampler, and the AMR. As mentioned before in response to another of the referee’s comments, we had neglected to share the specs of the machine that these times were recorded on, and this has now been added.

The reason for analysing the error as a function of the number of the likelihood evaluations, is that this method is intended for problems with expensive likelihoods, for which these overhead costs would be dwarfed. We have added a discussion on this to the beginning of Section 7.

Additionally to this, in Section 7.5, as detailed above, we have compared the CPU time that it takes the PAIS to reach a relative error in the mean of the posterior of 10^{-5} against that for trivially parallelised (non-communicative) RWMH. Even for a simple example such as this, PAIS outperforms standard parallelised MCMC by two orders of magnitude.

Referee 2 comments

Comment: *This work is an interesting contribution in the field of Monte Carlo algorithms. The paper is technically sound and well-written. Moreover, it is possible to note the effort and the care devoted by the authors..*

Answer: We thank the referee for their kind comments.

Comment: *However, I have some concerns.*

1) The starting-point algorithm in Table 1, is very similar to the basic scheme suggested in, for instance,

V. Elvira, L. Martino, D. Luengo, M. F. Bugallo, Improving Population Monte Carlo: Alternative Weighting and Resampling Schemes, Signal Processing Vol. 131, pp. 77-91, 2017.

In practice, you use a deterministic approach for the IS weights. I believe that the authors have analyzed and studied this idea independently, but this similarity should be discussed in the introduction, at least. You can also point out that you also propose some different resampling or adaptation schemes that are not contained in the previous work..

Answer: We agree that there are similarities to this other work, and we thank the referee for their comment. Additional references have been added, along with more detailed discussions about the novel aspects of this manuscript in the introduction and discussion.

Comment: *2) I do not understand why you use the word "parallel" in the name of the algorithm. In Algorithm 1, it is difficult to recognize the parallelization. The denominator in the weight that you employ, in fact, is a bottleneck for a possible parallelization. Please clarify this point (even in the text)..*

Answer: We have added clarifications on this point to Section 3. The referee should note that in our convergence plots, the reason that we compare the algorithms with respect to the number of likelihood evaluations is that we envisaged that this method would be of most use when the likelihood is expensive, for example when it requires a high resolution approximation of the solution of a PDE, and/or there is a very large amount of data to be taken into account. In this context, it would be useful to spread the work of computing the likelihood for the ensemble members across a cluster, hence our use of the word "parallel". We acknowledge and thank the referee for indicating the issues with the denominator. In the case where the likelihood is expensive, these overhead costs, along with that of the resampler, should no longer be the bottleneck. This was not adequately explained in the text, and as such we have added a further discussion in Section 3.

Comment: 3) Please, add a table of acronyms that can help the reader. Now it is a bit confusing..

Answer: We have added a glossary of acronyms to an appendix.

Comment: 4) The references should be completed with other important contributions (considering also the same kind of weights, in some cases):

- O.Cappe, R.Douc, A.Guillin, J.M.Marin, C.P.Robert, Adaptive importance sampling in general mixture classes, *Statistics and Computing* 18 (2008) 447-459. - R. Douc, A. Guillin, J. M. Marin, C. P. Robert, Convergence of adaptive mixtures of importance sampling schemes, *Annals of Statistics* 35 (2007) 420-448. - R. Douc, A. Guillin, J. M. Marin, C. P. Robert, Minimum variance importance sampling via population Monte Carlo, *ESAIM: Probability and Statistics* 11 (2007) 427-447. - L. Martino, V. Elvira, D. Luengo, J. Corander, Layered Adaptive Importance Sampling, *Statistics and Computing*, Vol. 27 (3), pp. 599-623, 2017..

Answer: We thank the referee for thier suggestions for other highly relevant references.

Associate Editor comments

Comment: A very similar version of the authors PAIS algorithm has been studied before, as explained by the second referee. Moreover, I could argue that the PAIS algorithm is a particular case since weights are considered deterministic. This is a serious issue in the paper, let alone that the authors are unaware of the series of papers mentioned by referee 2, that are closely related to their work. This should be explained in detail and the paper original perspective and contributions be further highlighted..

Answer: Following suggestions by the second referee and the associate editor, we have added a significant new section to the introduction to highlight the other ensemble importance schemes already in the literature. We have also highlighted further the contribution made by this manuscript in the introduction and discussion sections.

Comment: An other serious concern is why not use PMC-MCMC for the updating scheme in Algorithm 1. That would be as costly (since it also involves evaluating the posterior) but the y_i s will have the correct distribution..

Answer: This comment is not clear to us - what is the “updating scheme” in Algorithm 1? Also, the (weighted) y_i are distributed according to the posterior, as discussed in Section 6. We’d be happy to address this point if the AE is able to clarify it for us.

Comment: *The automatic tuning will be as useful as the what traditional adaptive tuning is: in the limit, the proposal cov matrix converges to the cov matrix of the target. This is ok for Gaussian type targets but will fail for multimodal, multiscale distributions, mostly present in non-trivial forward maps, as those suggested in section 2. Something similar happens with gradient based proposals, since from the onset they are not conceived to scape local maxima..*

Answer: This is an interesting point, and the subject of our most recent paper in this area, in which transport maps are used to stabilise and accelerate the PAIS methodology, and also addressing exactly the point raised by the AE. We have added a discussion on this topic in the discussion section, proposing it as an area for future study (our new manuscript will be completed in the next few weeks).

Comment: *In Algorithm 1, note that in step 5 the actual number of different particles could be drastically reduced and in some cases collapsed to 1. As mentioned by the authors, this is a basic concern in SIR type algorithms, that they may become totally inefficient. Then, why this seems not to bother the authors and no provision was included in their algorithm to mend this problem?.*

Answer: Provision is made within the code for such a catastrophic collapse of the weights, and it was remiss of us not to include details on this. We have added details on this to the end of Section 3.

Comment: *The Radon-Nokodyn version of Bayes theorem in (2.3) is unnecessary in the actual finite dimensional context the authors work with. The authors should state that, although there is an infinite dimensional version for the posterior distribution, their discussion is restricted to the parametric, finite dimensional case..*

Answer: We thank the AE for this comment. There is no need for this version of Bayes' theorem here, and so we have switched to the finite dimensional version.

Comment: *Note that in (2.4) the total measure of the space for \mathbf{z} is M (!)..*

Answer: Thanks to the AE for picking up on this typo, this has now been fixed.

PARALLEL ADAPTIVE IMPORTANCE SAMPLING

COLIN COTTER*, SIMON COTTER†, AND PAUL RUSSELL‡

Abstract. Markov chain Monte Carlo methods are a powerful and commonly used family of numerical methods for sampling from complex probability distributions. As applications of these methods increase in size and complexity, the need for efficient methods which can exploit the parallel architectures which are prevalent in high performance computing increases. In this paper, we aim to develop a framework for scalable parallel sampling algorithms. At each iteration, an importance sampling proposal distribution is formed using the ~~an~~ ensemble of particles. A stratified sample is taken from this distribution and weighted under the posterior, a state-of-the-art resampling method is then used to create an evenly weighted sample ready for the next iteration. We demonstrate that this parallel adaptive importance sampling (PAIS) method outperforms naive parallelisation of serial MCMC methods for low dimensional problems, and in fact shows better than linear improvements in convergence rates with respect to the number of ensemble members. We also introduce a new resampling strategy, approximate multinomial resampling (AMR), which while not as accurate as other schemes is substantially less costly for large ensemble sizes, which can then be used in conjunction with PAIS for complex problems. We also focus on how algorithmic parameters regarding the mixture proposal can be quickly tuned to optimise performance. In particular, we demonstrate this methodology's superior sampling for multimodal problems, such as those arising from inference for mixture models, and for problems with expensive likelihoods requiring the solution of a differential equation.

Key words. MCMC, parallel, importance sampling, Bayesian, inverse problems.

1. Introduction. Having first been developed in the early 1970s [23], Markov chain Monte Carlo (MCMC) methods have been of increasing importance and interest in the last 20 years or so. They allow us to sample from complex probability distributions which we would not be able to sample from directly. In particular, these methods have revolutionised the way in which inverse problems can be tackled, allowing full posterior sampling when using a Bayesian framework.

However, this often comes at a very high cost, with a very large number of iterations required in order for the empirical approximation of the posterior to be considered good enough. As the cost of computing likelihoods can be extremely large, this means that many problems of interest are simply computationally intractable.

This problem has been tackled in a variety of different ways. One approach is to construct increasingly complex MCMC methods which are able to use the structure of the posterior to make more intelligent proposals, leading to more thorough exploration of the posterior with fewer iterations. For example, the Hamiltonian or Hybrid Monte Carlo (HMC) algorithm uses gradient information and symplectic integrators in order to make very large moves in state with relatively high acceptance ~~probability~~probabilities [45]. Non-reversible methods are also becoming quite popular as they can improve mixing [3]. Riemann manifold Monte Carlo methods exploit the Riemann geometry of the parameter space, and are able to take advantage of the local structure of the target density to produce more efficient MCMC proposals [21]. This methodology has been successfully applied to MALA-type proposals and methods which exploit even higher order gradient information [6].

Since the clock speed of an individual processor is no longer following Moore's law [33], improvements in computational power are largely coming from the parallelisation of

*Department of Mathematics, Imperial College, London, UK

†School of Mathematics, University of Manchester, Manchester, UK. e: simon.cotter@manchester.ac.uk. SLC is grateful for EPSRC First grant award EP/L023393/1

‡School of Mathematics, University of Manchester, Manchester, UK

multiple cores. As such, the area of parallel MCMC methods is becoming increasingly of interest. One class of parallel MCMC method uses multiple proposals, with only one of these proposals being accepted. Examples of this approach include multiple try MCMC [28] and ensemble MCMC [34]. In [7], a general construction for the parallelisation of MCMC methods was presented, which demonstrated speed ups of up to two orders of magnitude when compared with serial methods on a single core. Another approach involves pre-fetching, where the possible future acceptances/rejections are calculated in advance [2].

A variety of other methods have been designed with particular scenarios in mind. For instance, sampling from high/infinite-dimensional posterior distributions is of interest in many applications. The majority of Metropolis-Hastings algorithms suffer from the curse of dimensionality, requiring more samples for a given degree of accuracy as the parameter dimension is increased. However some dimension-independent methods have been developed, based on Crank-Nicolson discretisations of certain stochastic differential equations [15]. Other ideas such as chain adaptation [22] and early rejection of samples can also aid reduction of the computational workload [47].

However, high dimensionality is not the only challenge that we may face. Complex structure in low dimensions can cause significant issues. These issues may arise due to large correlations between parameters in the posterior, leading to long thin curved structures which many standard methods can struggle with. These features are common, for example, in inverse problems related to epidemiology and other biological applications [24].

Multimodality of the posterior can also lead to incredibly slow convergence in many methods. Many methods allow for good exploration of the current mode, but the waiting time to the next switch of the chain to another mode may be large. Since many switches are required in order for the correct weighting to be given to each mode, and for all of the modes to be explored fully, this presents a significant challenge.

One application where this is an ever-present problem is that of mixture models. Given a dataset, where we know that the data is from two or more different distributions, we wish to be able to identify the parameters, e.g. the mean and variance and relative weight, of each part of the mixture [30]. The resulting posterior distribution is invariably a multimodal distribution, since the likelihood is invariant to permutations. Metropolis-Hastings algorithms, for example, will often fail to converge in a reasonable time frame for problems such as this. Since the posterior may be multimodal, independent of this label switching, it is important to be able to efficiently sample from the whole posterior.

Importance samplers are another class of methods which allow us to sample from complex probability distributions. A related class of algorithms, adaptive importance sampling (AIS) [27] reviewed in [5], had received less attention until their practical applicability was demonstrated in the mid-2000s [4, 9, 11, 25]. AIS methods produce a sequence of approximating distributions, constructed from mixtures of standard distributions, from which samples can be easily drawn. At each iteration the samples are weighted, often using standard importance sampling methods. The weighted samples are used to train the adapting sequence of distributions so that samples are drawn more efficiently as the iterations progress. The weighted samples form a sample from the posterior distribution under some mild conditions [31, 38].

Ensemble importance sampling schemes also exist, e.g. population Monte Carlo (PMC) [10]. PMC uses an ensemble to build a mixture or kernel density estimate (KDE) of the posterior distribution. The efficiency of this optimisation is restricted

by the component kernel(s) chosen, and the quality of the current sample from the posterior [8, 16, 17]. Extensions, such as [layered adaptive importance sampling](#) [32], adaptive multiple importance sampling algorithm (AMIS) [12], and adaptive population importance sampling (APIS) [31] have enabled these methods to be applied to various applications, including population genetics [46].

In this paper, we present a framework for parallelisation of importance sampling, which can be built around many of the current Metropolis-based methodologies in order to create an efficient target distribution from the current ensemble. The method makes use of a resampler based on optimal transport which has been used in the context of particle filters [37]. [We also detail how algorithmic parameters can be quickly tuned into optimal regimes, with respect to the effective sample size statistic.](#) In particular we demonstrate the advantages of this method when attempting to sample from multimodal posterior distributions, such as those arising from inference for mixture models. [The method that we present is specifically designed to tackle the challenges of certain types of low-dimensional inverse problem which have complex posterior structure, for example multimodality or strong correlations, and expensive likelihood evaluations. We will demonstrate that this methodology allows us to sample from such a posterior distribution with greater precision per likelihood evaluation than an uncorrelated ensemble of standard MCMC chains of the same type.](#)

[The method considered is similar to population Monte Carlo methods which have been previously studied, for example in \[19\]. However our approach exploits the current state-of-the-art resamplers. Higher accuracy resamplers, that provide a deterministic resample which optimally reflects the statistics of the original sample, lead to better mixture approximations of the posterior. These approximations can then be used to construct stable and efficient importance sampling proposals. The cost of the state-of-art resamplers can be prohibitive for large ensemble sizes, which are necessary for sampling in moderately higher dimensions, and as such we also present a greedy approximation of the optimal transport resampler, which provides us with good results for a fraction of the cost. We also detail how scaling parameters in the MCMC proposals can be quickly tuned, leading us to an efficient and fully-automated algorithm. We will demonstrate that despite additional overheads, the PAIS methodology can outperform its plain Metropolis-Hastings cousin by orders of magnitude to reach a histogram of a particular degree of convergence.](#)

In Section 2 we introduce some mathematical preliminaries upon which we will later rely. In Section 3 we present the general framework of the PAIS algorithm. In Section 4 we consider adaptive versions of PAIS which automatically tune algorithmic parameters concerned with the proposal distributions. In Section 5 we introduce the approximate multinomial resampling (AMR) algorithm which is a less accurate but faster alternative to resamplers which solve the optimal transport problem exactly. In Section 6 we consider consistency of the PAIS algorithm. In Section 7 we present some numerical examples, before a brief conclusion and discussion in Section 8.

2. Preliminaries. In this Section we will introduce preliminary topics and algorithms that will be referred to throughout the paper.

2.1. Bayesian inverse problems. In this paper, we focus on the use of MCMC methods for characterising posterior probability distributions arising from Bayesian inverse problems. We wish to learn about a particular unknown quantity x , of which we are able to make direct or indirect noisy observations. For now we say that x is a member of a Hilbert space X .

The quantity x is mapped on to observable space by the observation operator $\mathcal{G} : X \rightarrow \mathbb{R}^d$. We assume that the observations, D , are subject to Gaussian noise,

$$D = \mathcal{G}(x) + \varepsilon, \quad \varepsilon \sim \mu_\varepsilon = \mathcal{N}(0, \Sigma). \quad (2.1)$$

~~For example, if x is an initial condition for the wave equation, \mathcal{G} might calculate the height of a wave at a particular point and time.~~

These modelling assumptions allow us to construct the likelihood of observing the data D given the quantity $x = x^*$. Rearranging (2.1) and using the distribution of ε , we get:

$$\mathbb{P}(D|x = x^*) \propto \exp\left(-\frac{1}{2}\|\mathcal{G}(x^*) - D\|_\Sigma^2\right) = \exp(-\Phi(x^*)), \quad (2.2)$$

where $\|y_1 - y_2\|_\Sigma = (y_1 - y_2)^\top \Sigma^{-1} (y_1 - y_2)$ for $y_1, y_2 \in \mathbb{R}^d$.

~~As discussed in [14, 48], in order for this inverse problem to be well-posed in the Bayesian sense, we require μ_Y to be absolutely continuous with respect to the prior, μ_0 . A minimal regularity prior can be chosen informed by regularity results of the observational operator \mathcal{G} . Given such a prior, then the Radon-Nikodym derivative of the posterior measure, μ_Y , with respect to the prior measure, μ_0 , is proportional to the likelihood:~~

$$\frac{d\mu_Y}{d\mu_0} \propto \exp(-\Phi(x^*)).$$

that we are interested in characterising is given, up to a constant of proportionality, by

$$\pi_Y \propto \exp(-\Phi(x^*)) \pi_0(x^*),$$

where π_0 is the prior density on the quantity of interest x^* .

2.2. Particle filters and resamplers. In this subsection we briefly review particle filters, since the development of the resampler that we incorporate into the PAIS is motivated by this area.

Particle filters are a class of Monte Carlo algorithms designed to solve the filtering problem. That is, to find the best estimate of the true state of a system when given only noisy observations of the system. The solution of this problem has been of importance since the middle of the 20th century in fields such as molecular biology, computational physics and signal processing. In recent years the data assimilation community has contributed several efficient particle filters, including the ensemble Kalman filter (EnKF) [20] and the ensemble transform particle filter (ETPF) [37].

The ETPF defines a coupling T between two random variables Y and X , allowing us to use the induced map as a resampler. An *optimal coupling* T^* is one which maximises the correlation between X and Y [13]. This coupling is the solution to a linear programming problem in M^2 variables with $2M - 1$ constraints, where M is the size of the sample. Maximising the correlation preserves the statistics of X in the new sample.

In this work we use the resampler used within these particle filters ~~as resampling methods~~. We approximate the posterior distribution, μ_Y , with a ~~small~~ sample of

weighted particles, $\{(w_i, y_i)\}_{i=1}^M$. In filtering problems the weights would be found by incorporating new observed data whereas here we simply use importance weights. This distribution, $\hat{\mu}_Y \approx \mu_Y$, can then be resampled, using the ETPF or otherwise, into a new distribution, $\hat{\eta} \approx \mu_Y$,

$$\hat{\mu}_Y = \sum_{i=1}^M w_i \delta_{y_i}(\cdot) \xrightarrow{\text{ETPF}} \sum_{i=1}^M \frac{n_i}{M} \delta_{x_i}(\cdot) = \hat{\eta}. \quad (2.3)$$

where $n_i \in \{0, 1, \dots, M\}$ and $\sum_{i=1}^M n_i = M$. This new sample $\{x_i\}_{i=1}^M$ is equally weighted.

~~Particle~~ The resampler used within particle filters such as the ETPF are well suited to this problem since it is easy to introduce conditions in the resampling to ensure you obtain the behaviour you require. One downside is that the required ensemble size increases quickly with dimension, making it difficult to use in high-dimensional problems.

2.3. Deficiencies of Metropolis-type MCMC schemes. All MCMC methods are naively parallelisable. One can take a method and simply implement it ~~simultaneously~~ multiple times over an ensemble of processors. All of the states of all of the ensemble member can be recorded, and in the time that it takes one MCMC chain to draw N samples, M ensemble members can draw NM samples.

However, we argue that this is not an optimal scenario. First of all, unless we have a lot of information about the posterior, we will initialise the algorithm's initial state in the tails of the distribution. The samples that are initially made as the algorithm finds its way to the mode(s) of the distribution cannot be considered to be samples from the target distribution, and must be thrown away. This process is known as the burn-in. In a naively parallelised scenario, each ensemble member must perform this process independently, and therefore mass parallelisation makes no inroads to cutting this cost.

Moreover, many MCMC algorithms suffer from poor mixing, especially in multimodal systems. The ~~amount~~ number of samples that it takes for an MCMC trajectory to switch between modes can be large, and given that a large number of switches are required before we have a good idea of the relative probability densities of these different regions, it can be prohibitively expensive.

Another aspect of Metropolis-type samplers is that information computed about a proposed state is simply lost if we choose to reject that proposal in the Metropolis step. An advantage of importance samplers is that no evaluations of \mathcal{G} are ever wasted since all samples are saved along with their relative weighting.

~~Moreover, a naively parallelised MCMC scheme is exactly that -- naive. Intuition suggests that we can gain some speed-up by sharing information across the ensemble members, and this is what we wish to demonstrate in this paper.~~

These deficiencies of the naive method of parallelising MCMC methods motivated the development of the Parallel Adaptive Importance Sampler (PAIS). In the next section we will introduce the method in its most general form.

3. The Parallel Adaptive Importance Sampler (PAIS). Importance sampling can be a very efficient method for sampling from a probability distribution. A proposal density is chosen, from which we can draw samples. Each sample is assigned a weight given by the ratio of the target density and the proposal density at that point. They are efficient when the proposal density is concentrated in similar areas

to the target density, and incredibly inefficient when this is not the case. The aim of the PAIS is to use an ensemble of states to construct a proposal distribution which will be as close as possible to the target density. If this ensemble is large enough, the distribution of states will be approximately representative of the target density.

The proposal distribution could be constructed in many different ways, but we choose to use a mixture distribution, made up of a sum of MCMC proposal kernels. ~~This corresponds to replacing the Dirac distribution in the right-hand side of Equation 2.3 with a chosen distribution~~ e.g. Gaussian Gaussians centred at the current state, as in RWMH. Once the proposal is constructed, we can sample a new ensemble from the proposal distribution, and each is assigned a weight given by the ratio of the target density and the proposal mixture density. Assuming that our proposal distribution is a good one, then the variance of the weights will be small, and we will have many useful samples. Finally, we need to create a set of evenly weighted samples which best represent this set of weighted samples. This is achieved by implementing a resampling algorithm. These samples are not stored in order to characterise the posterior density, since the resampling process inevitably adds some error/~~bias~~. They are simply needed in order to inform the mixture proposal distribution for the next iteration of the algorithm.

Initially we will use the ETPF resampler algorithm [37], although we will suggest an alternative strategy in Section 5 where a large ensemble is required, and for which the ETPF may become more expensive. The resampling algorithm gives us a set of evenly weighted samples which represents the target distribution well, which we can use to iterate the process again. The algorithm is summarised in Algorithm 1.

We wish to sample states $x \in X$ from a posterior probability distribution μ_D , where D represents our data which we wish to use to infer x . Since we have M ensemble members, we represent the current state of all of the Markov chains as a vector $\mathbf{X} = [x_1, x_2, \dots, x_M]^\top$. We are also given a transition kernel $\nu(\cdot, \cdot)$, which might come from an MCMC method, for example the random walk Metropolis-Hastings proposal density $\nu(\cdot, x) \sim \mathcal{N}(x, \beta^2)$, where $\beta^2 \in \mathbb{R}$ defines the variance of the proposal.

Algorithm 1: The parallel adaptive importance sampler (PAIS).

- 1 Set $\mathbf{X}^{(0)} = \mathbf{X}_0 = [x_1^{(0)}, x_2^{(0)}, \dots, x_M^{(0)}]^\top$.
 - 2 **for** $i = 1, \dots, N$ **do**
 - 3 Sample $\mathbf{Y}^{(i)} = [y_1^{(i)}, y_2^{(i)}, \dots, y_M^{(i)}]^\top$, $y_j^{(i)} \sim \nu(\cdot; x_j^{(i-1)})$.
 - 4 Calculate $\mathbf{W}^{(i)} = [w_1^{(i)}, w_2^{(i)}, \dots, w_M^{(i)}]^\top$, $w_j^{(i)} = \frac{\pi(y_j^{(i)})}{\chi(y_j^{(i)}; \mathbf{X}^{(i-1)})}$, where

$$\chi(\cdot; \mathbf{X}^{(i-1)}) = \frac{1}{M} \sum_{j=1}^M \nu(\cdot; x_j^{(i-1)}).$$
 - 5 Resample: $(\mathbf{W}^{(i)}, \mathbf{Y}^{(i)}) \rightarrow (\frac{1}{M}\mathbf{1}, \mathbf{X}^{(i)})$.
 - 6 Output (\mathbf{W}, \mathbf{Y}) .
-

Since the resampling does not give us a statistically identical sample to that which is ~~inputted input~~, we cannot assume that the samples $\mathbf{X}^{(i)}$ are samples from the posterior. Therefore, as with serial importance samplers, the weighted samples ~~$(\mathbf{W}^{(i)}, \mathbf{Y}^{(i)})_{i=1}^N$~~ $(\mathbf{W}, \mathbf{Y})_{i=1}^N$ are the samples from the posterior that we will analyse.

In each iteration of the algorithm, we are required to compute the likelihood for each member of the ensemble. In the case where the likelihood is expensive to compute, for example because it requires the numerical approximation of an ODE or PDE, and/or there is a very large amount of observations in the data, it will be expedient to split this computational effort across multiple cores. Whether this is efficient will depend on the architecture of the machine, and the relative cost of communication across the cores, to the cost of the likelihood evaluation. Evaluation of the denominator in the weights also contributes to the overhead costs of this approach over and above a standard MCMC method, but this too can be parallelised efficiently, depending on communication costs of the architecture of the machine, reducing the $\mathcal{O}(M^2)$ cost to $\mathcal{O}(M)$ over M cores, for example.

The key is to choose a suitable transition kernel ν such that if $\mathbf{X}^{(i)}$ is a representative sample of the posterior, then the mixture density $\chi(\cdot; \mathbf{X}^{(i)})$ is a good approximation of the posterior distribution. If this is the case, the newly proposed states $\mathbf{Y}^{(i)}$ will also be a good sample of the posterior with low variance in the weights $\mathbf{W}^{(i)}$.

In Section 7, we will demonstrate how the algorithm performs, primarily using random walk (RW) proposals. We do not claim that this choice is optimal, but is simply chosen as an example to show that sharing information across ensemble members can improve on the original MH algorithm and lead to tolerances being achieved in fewer evaluations of \mathcal{G} . This is important since if the inverse problem being tackled involves computing the likelihood from a very large data set, or where the likelihood requires the numerical solution of a differential equation, this could lead to a large saving of computational cost. We have observed that using more complex (and hence more expensive) kernels ν , does not significantly improve the speed of convergence of the algorithm for posteriors that we have considered, although kernels such as those used in MALA can be more stable for certain problems [44].

Care needs to be taken when choosing kernel(s) for the proposal distribution to ensure that the proposals are overall mixture is absolutely continuous with respect to the posterior distribution. If this is not the case, importance weights grow exponentially for samples in the tails can become high variance where samples are made away from modes in the distribution. This will drastically decrease the efficiency of sampling.

It is also possible, especially during the early stages of the algorithm, where the mixture proposal is a poor approximation of the target distribution, for a sample to be produced with a weight which is orders of magnitude bigger than the rest. This is a problem, since the resampling step will then lead to a vector of particles all centred at the outlier sample. Fortunately these spikes in weights can easily be detected, and in such a case, the sample for this iteration can be removed. The ETPF can then be used to formulate a new mixture sample with one particle placed at the outlier point, and the others distributed according to the particle positions before the problematic iteration. This ensures that the region where the problem occurred is better represented by the mixture, and prevents further spikes in this region.

4. Automated tuning of algorithm parameters. Efficient selection of scaling parameters in MCMC algorithms is critical to achieving optimal mixing rates and hence achieving fast convergence to the target density. It is well known that a scaling parameter which is either too large or too small results in a Markov chain with high autocorrelation. One aspect worthy of consideration with the PAIS, is finding an appropriate proposal kernel ν such that the mixture distribution χ is a close approximation to the posterior density π .

Most MCMC proposals have parametric dependence which allows the user to control their variance. For example, in the RW proposal $y = x + \beta\eta$, the parameter β is the standard deviation of the proposal distribution. Therefore the proposal distributions can be tuned such that they are ~~are~~ slightly over-dispersed. This tuning can take place during the burn-in phase of the algorithm. Algorithms which use this method to find optimal proposal distributions are known as adaptive MCMC algorithms, and have been shown to be convergent provided that they satisfy certain conditions [41, 42]. Alternatively, a costly trial and error scheme with short runs of the MCMC algorithm can be used to find an acceptable value of β .

Algorithms which use mixture proposals, e.g. PAIS, must tune the variance of the individual kernels within the proposal mixture. This adaptivity during early iterations has some added benefits over and above finding an optimal parameter regime for the algorithm. If the initial value of the proposal variance is chosen to be very large, then the early mode-finding stages of the algorithm are expedited. Adaptively reducing the proposal variances to an optimal value then allows us to explore each region efficiently. Using an ensemble of chains allows quick and effective assessment of the value of the optimal scaling parameter.

In many MCMC algorithms such as the Random Walk Metropolis-Hastings (RWMH) algorithm, the optimal scaling parameter can be found by searching for the parameter value which gives an optimal acceptance rate, e.g. for near Gaussian targets the optimal rates are 23.4% for RWMH and 57.4% for MALA [40]. This method is not applicable to PAIS so we must use other statistics to optimise the scaling parameter. Section 4.1 gives some possible methods for tuning β .

4.1. Statistics for Determining the Optimal Scaling Parameter.

4.1.1. Determining optimal scaling parameter using error analysis. When available, an analytic form for the target distribution allows us to assess the convergence of sampling algorithms to the target distribution. Common metrics for this task include the relative error between the sample moments and the target's moments, or the relative L^2 error between the sample histogram and the target density, $\pi(x|D)$. The relative error in the n -th moment, \hat{m}_n , is given by:

$$e = \left| \frac{\hat{m}_n - \mathbb{E}[X^n]}{\mathbb{E}[X^n]} \right|, \quad \text{where} \quad \hat{m}_n = \frac{1}{N} \sum_{i=1}^N x_i^n, \quad (4.1)$$

and $\{x_i\}_{i=1}^N$ is a sample of size N .

The relative L^2 error, E , between a continuous density function to a piecewise constant approximation of that density, can be defined by considering the difference in mass between the self-normalised histogram of the samples and the posterior distribution over a set of disjoint sets or "bins":

$$E^2 = \sum_{i=1}^{n_b} \left[\int_{R_i} \pi(s|D) ds - vB_i \right]^2 / \sum_{i=1}^{n_b} \left[\int_{R_i} \pi(s|D) ds \right]^2, \quad (4.2)$$

where the regions $\{R_i\}_{i=1}^{n_b}$ are the d -dimensional histogram bins, so that $\bigcup_i R_i \subseteq X$ and $R_i \cap R_j = \emptyset$, n_b is the number of bins, v is the volume of each bin, and B_i is the value of the i th bin. This metric converges to the standard definition of the relative L^2 error as $v \rightarrow 0$.

These statistics cannot be used in general to find optimal values of β since they require the analytic solution, and are expensive to approximate. However they can be used

to assess the ability of other indicators to find the optimal scaling parameters in a controlled setting.

4.1.2. The effective sample size. The effective sample size, n_{eff} , can be used to assess the efficiency of importance samplers. Ideally, in each iteration, we would like all M of our samples to provide us with new information about the posterior distribution. In [practisepractice](#), we cannot achieve a perfect effective sample size of M .

The effective sample size of a weighted sample is defined in the following way:

$$n_{\text{eff}} = \frac{\left(\sum_{i=1}^M w_i\right)^2}{\sum_{i=1}^M w_i^2} \approx \frac{M\mathbb{E}(w)^2}{\mathbb{E}(w^2)} = M \left(1 - \frac{\text{var}(w)}{\mathbb{E}(w^2)}\right).$$

The second two expressions are true when $M \rightarrow \infty$. From the last expression, if the variance of the weights is zero then $n_{\text{eff}} = M$; this is our ideal scenario. In the limit $M \rightarrow \infty$, maximising the effective sample size is equivalent to minimising the variance of the weights.

The statistic n_{eff} is easier to deal with than the variance of the weights, as it takes values on $[1, M]$ while the variance of the weights takes values on \mathbb{R}_{\geq} . Moreover the variance of the weights can vary over many orders of magnitude causing numerical [issuesinstabilities](#), so that the effective sample size is more desirable as an indicator of optimality.

In this paper we tune our simulations using the effective sample size statistic. We calculate this statistic using a sample size of Mn_k , where $1 \leq n_k \leq N$ is sufficiently large enough to obtain a reasonable estimate of n_{eff} with scaling parameter δ_k . Calculating n_{eff} over these subsets of the simulations tends to underestimate the optimal value of the scaling parameter due to the possibility of missing unlikely proposals with extreme weights.

The effective sample size also has another useful property; if we consider the algorithm in the early stages, for example, we have all M particles in the tails of the target searching for a mode. The particle closest to the mode will have an exponentially higher weight and the effective sample size will be close to 1. In later stages the ensemble populates high density regions, and better represents the posterior distribution. This leads to smaller variation in the weights, and a higher effective sample size. By looking for approximations of the effective sample size which look like a stationary distribution, we can tell when the algorithm is working efficiently.

Figure 4.1 demonstrates that the effective sample size flattens out as the relative L^2 difference between the posterior distribution and the proposal distribution stabilises close to its minimum.

Figure 4.2 shows how the effective sample size used in PAIS is affected by the ensemble size. We see from subfigure (a) that as the ensemble size increases, the optimal scaling parameter decreases. This is expected since the larger ensemble allows for finer resolution in the proposal distribution approximation of the posterior distribution. We also see that the algorithm becomes less sensitive to changes in the scaling parameter as the ensemble size increases. Subfigure (b) shows that as the ensemble size increases, the algorithm becomes more efficient.

4.2. Adaptively Tuned PAIS. To adapt the scaling parameter β , we use a version of the gradient ascent method modified for a stochastic function. Some more sophisticated examples are described in [1, 26, 42].

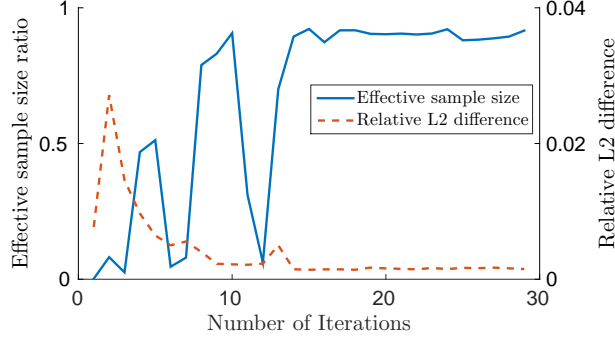
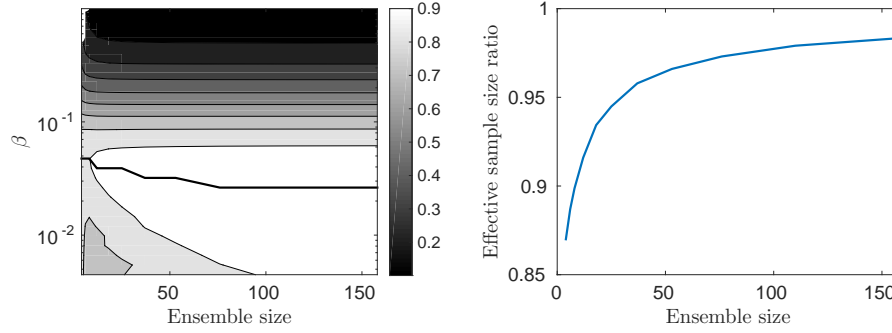


FIG. 4.1. The effective sample size ratio and relative L^2 difference E during the first 30 iterations. These numerics are taken from the example in Section 7.4.3 using the PAIS-RW algorithm.



(a) Contours showing optimal ranges of the scaling parameter. (b) The highest effective sample size achieved for each ensemble size.

FIG. 4.2. The behaviour of the effective sample size as the ensemble size increases, considering the example in Section 7.2 using the PAIS-RW algorithm.

Algorithm 2: Adaptively tuned PAIS algorithm.

- 1 Define update times $\{n_k\}_k$.
- 2 **for** $n = 1, \dots, N$ **do**
- 3 Complete steps 3-5 of Algorithm 1.
- 4 **if** $n \in \{n_k\}$ **then**
- 5 Divide ensemble into two halves. Use these halves to estimate the gradient in n_{eff} at β_k .
- 6 Update β_k using gradient ascent,

$$\beta_{k+1} = \beta_k + \gamma \nabla n_{\text{eff}}.$$

Our adaptive algorithm is given in Algorithm 2. We choose update times which, as suggested in Section 4.1.2 allow for a reasonable estimate of the effective sample size, but do not waste too many iterations. In Step 6, the gradient ascent parameter γ

may decrease over time, e.g. as a function of n_k .

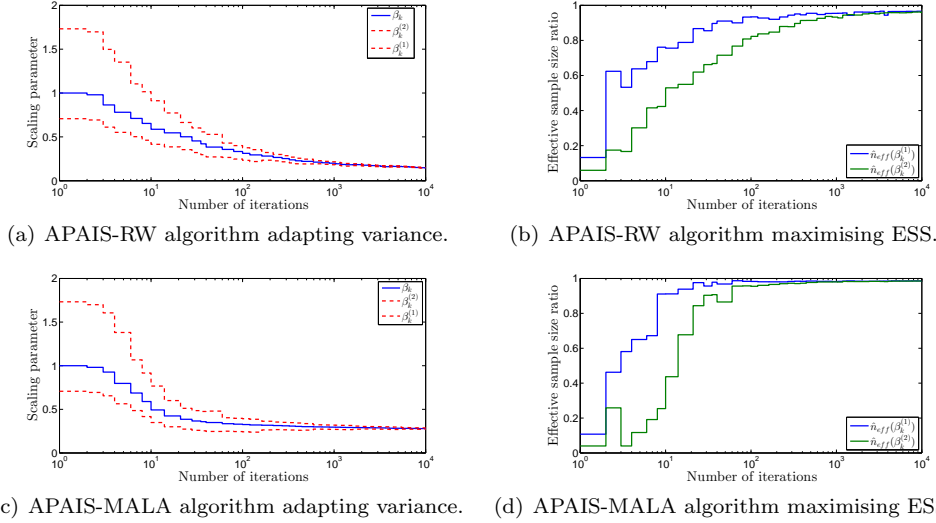


FIG. 4.3. *Demonstration of the convergence of the adaptive scaling parameter in the adaptively tuned PAIS algorithms.*

In the context of an expensive likelihood evaluation, slow convergence of the algorithmic parameters in this adaptive regime would be a concern. However, this is rarely an issue, due in part to the speed at which the PAIS algorithm completes the burn-in phase. Figure 4.3 demonstrates this for two different versions of the PAIS methodology with a RW proposal and a MALA [39] proposal. This picture is typical in our experience, with the scaling parameters converging in $\mathcal{O}(10^2)$ iterations to an optimal regime.

5. Approximate Multinomial Resampling. Although the ETPF is optimal in terms of preserving statistics of the sample, it produces the optimal linear coupling, it can also become quite costly as the number of ensemble members is increased. It is arguable that in the context of PAIS, we do not require this degree of accuracy, and that a faster more approximate method for resampling could be employed. One approach would be to use the bootstrap resampler, which simply takes the M ensemble members' weights and constructs a multinomial distribution, from which M samples are drawn. This is essentially the cheapest resampling algorithm that one could construct. However it too has some drawbacks. The algorithm is random, and as such it is possible for all of the ensemble members in a particular region not to be sampled. This could be particularly problematic when attempting to sample from a multimodal distribution, where it might take a long time to find one of the modes again. The bootstrap filter is also not guaranteed to preserve the mean of the weighted sample, unlike the ETPF.

Ideally, we would like to use a resampling algorithm which is not prohibitively costly for moderately or large sized ensembles, which preserves the mean of the samples, and which makes it much harder for the new samples to forget a significant region in the density. This motivates the following algorithm, which we refer to as **approximate multinomial resampling**—the approximate multinomial resampler (AMR), which is a greedy approximation of the ETPF resampler.

Instead of sampling M times from an M -dimensional multinomial distribution as is the case with the bootstrap algorithm, we sample once each from M different multinomials. Suppose that we have M samples y_n with weights w_n . The multinomial sampled from in the bootstrap filter has a vector of probabilities given by:

$$\frac{1}{\sum w_n} [w_1, w_2, \dots, w_M] = \bar{\mathbf{w}},$$

with associated states y_n . We wish to find M vectors $\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_M\} \subset \mathbb{R}_{\geq 0}^M$ such that $\frac{1}{M} \sum \mathbf{p}_i = \bar{\mathbf{w}}$. The AMR is then given by a sample from each of the multinomials defined by the vectors $\mathbf{p}_i = [p_{i,1}, p_{i,2}, \dots, p_{i,M}]$ with associated states \mathbf{y}_i . Alternatively, as with the ETPF, a deterministic sample can be chosen by picking each sample to be equal to the mean value of each of these multinomial distributions, i.e. each new sample \hat{x}_i is given by:

$$\hat{x}_i = \sum p_{i,j} x_j, \quad i \in \{1, 2, \dots, M\}. \quad (5.1)$$

The resulting sample has several properties which are advantageous in the context of being used with the PAIS algorithm. Firstly, we have effectively chopped up the multinomial distribution used in the bootstrap filter into M pieces, and we can guarantee that exactly one sample will be taken from each section. This leads to a much smaller chance of losing entire modes in the density, if each of the sub-multinomials is picked in an appropriate fashion. Secondly, if we do not make a random sample for each multinomial with probability vector \mathbf{p}_i but instead take the mean of the multinomial to be the sample, this algorithm preserves the mean of the sample exactly. Lastly, as we will see shortly, this algorithm is significantly less computationally intensive than the ETPF.

There are of course infinitely many different ways that one could use to split the original multinomial up into M parts, some of which will be far from optimal. The method that we have chosen is loosely based on the idea of optimal transport. We search out states with the largest weights, and choose a cluster around these points based on the closest states geographically. This method is not optimal since once most of the clusters have been selected the remaining states may be spread across the parameter space.

Algorithm 3: The approximate multinomial resampler (AMR).

```

1  $\mathbf{z} = M\bar{\mathbf{w}}.$ 
2 for  $i = 1, \dots, M$  do
3    $J = \arg \max_j z_j.$ 
4    $p_{i,J} = \min\{1, z_J\}.$ 
5    $z_J = z_J - p_{i,J}.$ 
6   while  $\sum_j p_{i,j} < 1$  do
7      $K = \arg \min_{k \in \{k | z_k > 0\}} \|y_J - y_k\|.$ 
8      $p_{i,K} = \min\{1 - \sum_j p_{i,j}, z_K\}.$ 
9      $z_K = z_K - p_{i,K}.$ 
10   $x_i = \sum_k p_{i,k} y_k.$ 
```

Algorithm 3 describes the basis of the algorithm with deterministic resampling, using the means of each of the sub-multinomials as the new samples. This resampler

was designed with the aims of being numerically cheaper than the ETPF, and more accurate than straight multinomial resampling. Therefore we now present numerical examples which demonstrate this.

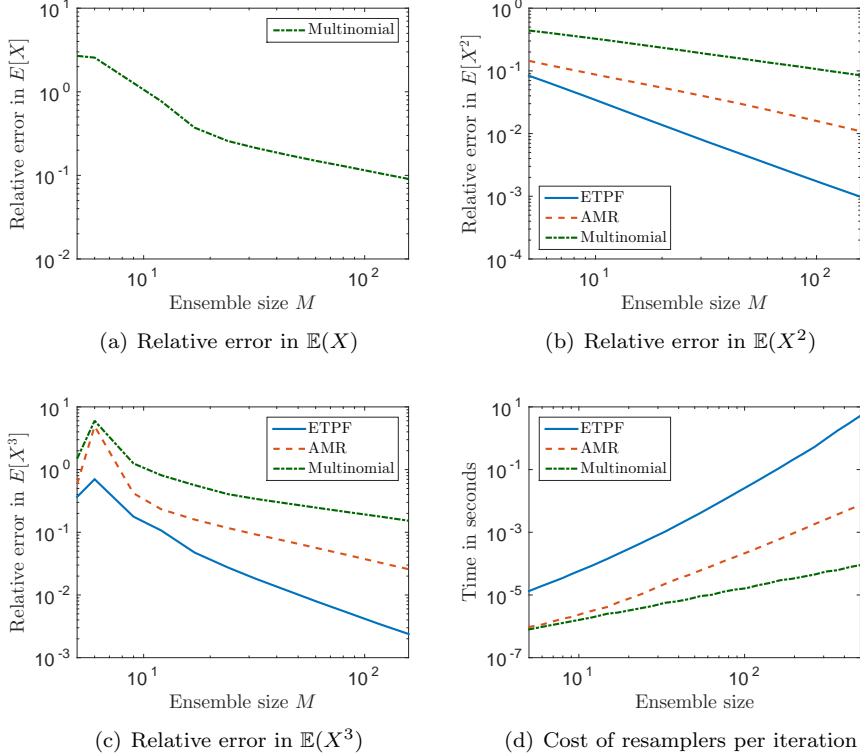


FIG. 5.1. Comparison of the performance between different resampling schemes. The example in Section 7.2 is implemented for this demonstration.

To test the accuracy and speed of the three resamplers (ETPF, bootstrap and AMR), we drew a sample of size M from the proposal distribution $\mathcal{N}(1, 2)$. Importance weights were assigned, based on a target distribution of $\mathcal{N}(2, 3)$. The statistics of the resampled outputs were compared with the original weighted samples. Figure 5.1 (a)-(c) show how the relative errors in the first three moments of the samples changes with ensemble size M for the three different ~~samplers~~resamplers. As expected, the AMR lies somewhere between the high accuracy of the ETPF and the less accurate bootstrap resampling. Note that only the error for the bootstrap multinomial sampler is presented for the first moment since both the ETPF and the AMR preserve the mean of the original weighted samples up to machine precision. Figure 5.1 (d) shows how the computational cost, measured in seconds, scales with the ensemble size for the three different methods, where timings have been taken from simulations on a Dell server with four 8-core 3.3GHz CPUs and 64Gb memory. These results demonstrate that the AMR behaves how we wish, and importantly ensures that exactly one sample of the output will lie in each region with weights up to $\frac{1}{M}$ of the total.

We will use the AMR in the numerics in Section 7.4 where we have chosen to use a larger ensemble size. We do not claim that the AMR is the optimal choice within PAIS, but it does have favourable features, and demonstrates how different choices of

resampler can affect the speed and ~~accuracy~~-efficiency of the PAIS algorithm.

6. Consistency of PAIS. As outlined in [31] consistency of population AIS algorithms can be considered in two different cases. In the first case we fix the number of iterations $N < \infty$, but allow the population size to grow to infinity $M \rightarrow \infty$. In the second case we hold the population size fixed $M < \infty$, and allow infinitely many iterations $N \rightarrow \infty$.

In case one, from standard importance sampling results we know that for an iteration n , as $M \rightarrow \infty$, we obtain a consistent estimator for any statistic of interest,

$$\hat{\phi}(\Theta) \approx \frac{1}{\hat{Z}} \sum_{i=1}^M \frac{1}{M} w_i \phi(\theta_i) \rightarrow \phi(\Theta),$$

where the normalisation constant, $\hat{Z} = \frac{1}{M} \sum_{i=1}^M w_i$, also converges to the true normalisation constant Z [38].

Case two is slightly more involved. Estimation of the normalisation constant Z is biased, and so estimates of statistics are sums of independent but biased estimators. Since the estimators are independent, the proof of consistency of PAIS in this second limit can be approached in the same way as the PMC algorithm, where it has been shown that $\hat{Z} \rightarrow Z$ [38]. Since the normalisation constant is consistent, sums of the independent estimators are also consistent.

7. Numerical Examples. The numerical examples in this section were all computed on a single core of a Dell server with four 8-core 3.3GHz CPUs and 64Gb memory. Throughout we measure the efficiency of the algorithms in terms of the number of iterations (i.e. likelihood evaluations) which are required in order to reach a given order of accuracy. The reason for this is that this method is designed with a particular type of challenging inverse problem in mind, namely one which is low dimensional, has an expensive likelihood (such as a PDE solve, and/or very large data) which dwarfs the overhead cost of PAIS, and which could have complex posterior structure, such as correlations, multimodality or other highly non-Gaussian features.

In this section we demonstrate convergence of the PAIS algorithm. We also investigate the convergence properties of PAIS compared with naively parallelised Metropolis-Hastings samplers. We assess the performance of the PAIS algorithm using the relative L^2 error defined in Equation (4.2), as well as the relative error in the first moment (Equation (4.1)).

7.1. Numerical implementation. For each example, we perform the following three tasks.

1. **Finding the optimal scaling parameters:** We chose 32 values of β evenly spaced on a log scale in the interval $[10^{-5}, 2]$. We ran the PAIS and MH algorithms for one million iterations, each with an ensemble size of $M = 50$. We took 32 repeats and used the geometric means of the sample statistics in Section 4.1 to find the optimal scaling parameter.
2. **Measuring convergence of nonadaptive algorithms:** We ran the PAIS and MH algorithms, using the scaling parameters found in Step 1, for one million iterations, again with $M = 50$. The simulations are repeated 32 times. The performance is evaluated using the relative L^2 error defined in Equation (4.2).

3. **Measuring convergence of adaptive algorithms:** We run the adaptive algorithms under the same conditions as the nonadaptive algorithms, and again use the relative L^2 error to compare efficiency. The adaptive algorithms are initialised with $\beta_1 = 1$.

7.2. Sampling from a one dimensional Gaussian distribution. This first example shows how PAIS handles searching for the mode of a Gaussian distribution when the initial state is a long way out in the tails. We optimise the naively parallelised RWMH algorithm using the optimal acceptance rate $\hat{\alpha} = 0.5$. This value differs from the theoretical asymptotic value of 0.234 which applies in higher dimensions. This higher acceptance rate is commonly used for one dimensional Gaussian posteriors [43]. To tune the PAIS algorithm we maximise the effective sample size as discussed in Section 4.1.2.

7.2.1. Target distribution. Consider the case of a linear observation operator which maps a one dimensional parameter onto the observable space, $\mathcal{G}: x \mapsto x$. We assign centred priors to the parameter x , as well as to the observational noise. We choose the true value of the parameter to be $x_{\text{ref}} = 4$, and choose $\sigma^2 = \tau^2 = 0.1$. Then following Equations (2.1) and (2.2),

$$\text{law}(\mu_D) = \pi(x|D) \propto \exp\left(-\frac{1}{2}\|x - D\|_{\sigma^2}^2 - \frac{1}{2}\|x\|_{\tau^2}^2\right). \quad (7.1)$$

This set-up results in a posterior density in which the vast majority of the density is out in the tails of the prior distribution. The Kullback-Leibler (KL) divergence, which gives us a measure of how different the prior and posterior are, is $D_{KL}(\mu_D||\mu_0) = 4.67$ for this problem. A KL divergence of zero indicates that two distributions are identical almost everywhere.

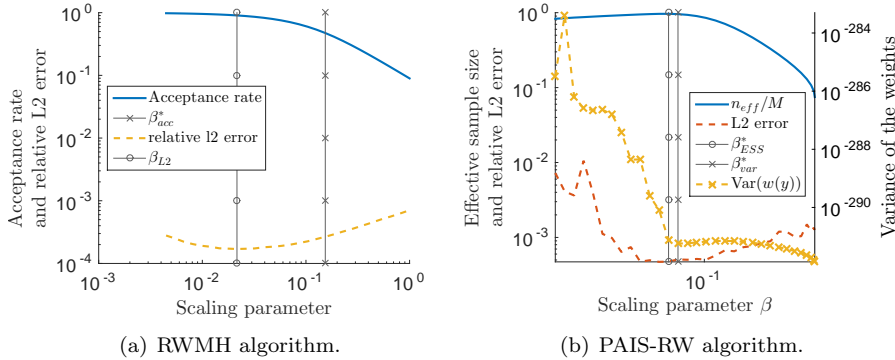


FIG. 7.1. Finding optimal values of β for the example in Section 7.2. The setup is as in Section 7.1. Resampling is performed using the ETPF.

7.2.2. Optimising the scaling parameter. Figure 7.1 (a) shows the two values of β which are optimal according to the acceptance rate and relative L^2 error criteria for the RWMH algorithm. The smaller estimate comes from the relative L^2 error, and the larger from the acceptance rate. The results in Figure 7.1 are summarised in Table 7.1. Since in general we cannot calculate the relative L^2 error, we must optimise the algorithm using the acceptance rate. From the relative L^2 error

Statistic	RWMH	Statistic	PAIS-RW
$\beta_{L^2}^*$	2.1e-2	β_{eff}^*	4.7e-2
$\beta_{\%}^*$	1.5e-1	$\beta_{\text{var}(w(y))}^*$	5.8e-2
Acceptance Rate ($\beta_{L^2}^*$)	9.0e-1	$\beta_{L^2}^*$	3.9e-2
Acceptance Rate ($\beta_{\%}^*$)	5.0e-1		

TABLE 7.1

Optimal values of β summarised from Figure 7.1. Statistics calculated as described in Section 4.1. The values $\beta_{L^2}^*$ and $\beta_{\%}^*$ are the optimal scaling parameters found by optimising the relative L^2 errors and acceptance rate respectively. Similarly β_{eff}^* and $\beta_{\text{var}(w(y))}^*$ optimise the effective sample size and variance of the weights statistics.

curve we can see that the minimum is very wide and despite the optimal values being very different there is not a large difference in the convergence rate.

Figure 7.1 (b) shows the effective sample size ratio compared to the error analysis and the variance of the weights. The relative L^2 error graph is noisy, but it is clear that the maximum in the effective sample size and the minimum in the variance of the weights are both close to the minimum in the relative L^2 error. Due to this we say that the estimate of the effective sample size found by averaging the statistic over each iteration is a good indicator for the optimal scaling parameter.

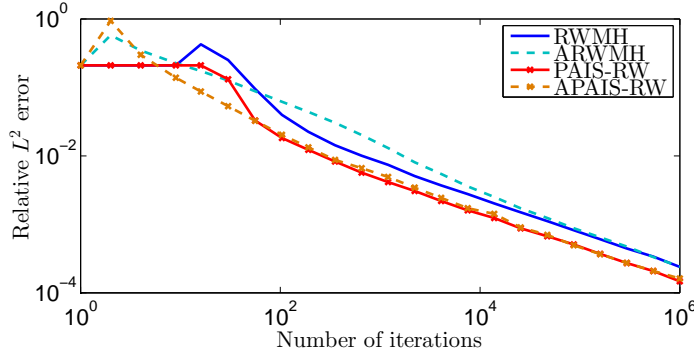


FIG. 7.2. Convergence of the (A)RWMH and (A)PAIS-RW algorithms. The setup is as in Section 7.1 (2, 3). Resampling is performed using the ETPF.

7.2.3. Convergence of RWMH vs PAIS-RW. Figure 7.2 shows that the PAIS-RW algorithm converges to the posterior distribution significantly faster than the RWMH algorithm, in both L^2 error and relative error in the moments. A description of the speed up attained by this algorithm is given in Section 7.3.4. The adaptive algorithms are also shown in Figure 7.2. We see that both adaptive algorithms converge to the posterior at a similar speed to the respective optimised algorithm. This shows that, particularly for PAIS, we can optimise simulations efficiently on the fly.

7.2.4. Scaling of the PAIS algorithm with ensemble size. In this and the next example we use an ensemble size $M = 50$, however it is interesting to see how the PAIS algorithm scales when we increase the ensemble size, and if there is a minimum required ensemble size. We implement the problem in Section 7.2, using the RWMH and PAIS-RW algorithms with ensemble sizes in the interval $M \in [1, 160]$.

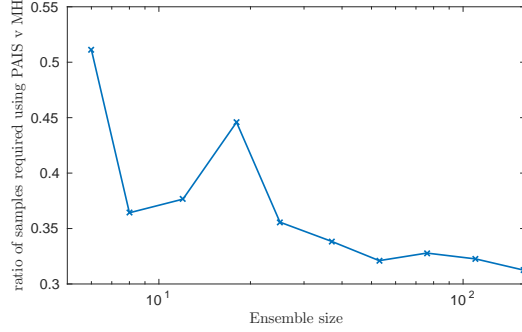


FIG. 7.3. Ratio of PAIS-RW samples required to reach the same tolerance as the RWMH algorithm, using the ETPF resampler.

Figure 7.3 was produced using the method of finding optimal β described in Section 7.1 (1), then running 32 repeats at each ensemble size

We repeat the process described in Section 7.1 for each of the ensemble sizes we are interested in. The convergence rates are then found by regressing through the data. The graph is still very noisy but demonstrates that increasing the ensemble size continues to reduce the number of iterations required in comparison with naively parallelised MH. The decreasing trend shows superlinear improvement of PAIS with respect to ensemble size, in terms of the number of iterations required, which is a demonstration of our belief that parallelism of MCMC should give us added value over and above that provided by naive parallelism. This decrease is due to the increasing effective sample size shown in Figure 4.2 (b).

7.3. Sampling from Bimodal Distributions. In this second example we investigate the behaviour of the PAIS algorithm when applied to bimodal problems. MH methods can struggle with multimodal problems, particularly where switches between the modes are rare, resulting in incorrectly proportioned modes in the histograms. This example demonstrates that the PAIS algorithm redistributes particles to new modes as they are found. This means that we expect the number of particles in a mode to be approximately proportional to the probability density in that mode, resulting in faster global convergence.

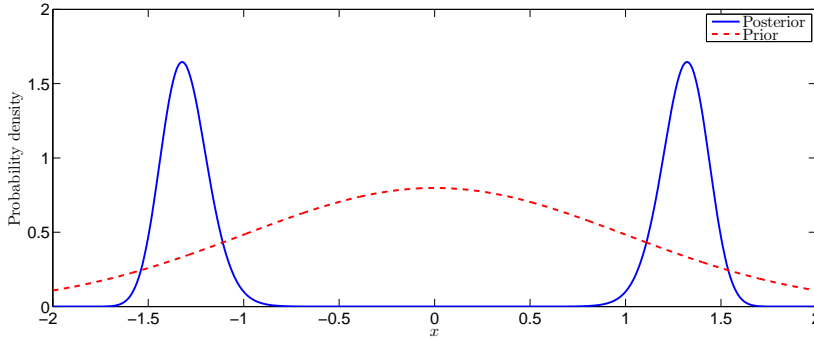


FIG. 7.4. Posterior distribution for the bimodal example in Section 7.3.1.

7.3.1. Target distribution. In this example we look at a bimodal posterior distribution which has a ~~high-energy barrier~~ region of very low probability density between the two modes. In the MH algorithm, the expected number of iterations for a switch between modes is a large proportion of the total number of iterations. This posterior is presented in Figure 7.4.

We consider a non-linear observation operator $\mathcal{G}: x \mapsto x^2$, and assign the prior $\mu_0 = \mathcal{N}(0, \tau^2 = 0.25)$. Again, we assume the true parameter $x_{\text{ref}} = 2$ is observed according to Equation (2.1) with observational noise $\sigma^2 = 0.1$, and the posterior is

$$\pi(x|D) \propto \exp\left(-\frac{1}{2\sigma^2}\|x^2 - D\|^2 - \frac{1}{2\tau^2}\|x\|^2\right).$$

7.3.2. Calculating values of the optimal scaling parameter. Calculating the optimal values of the scaling parameters for this example is similar to the previous example. These values are given in Table 7.2. The subscript on β refers to the criterion which has been optimised.

Algorithm	β_{acc}^*	β_{eff}^*	$\beta_{L^2}^*$ <u>$\beta_{L^2}^*$</u>
RWMH	2.3e-1	-	9.3e-1
PAIS-RW	-	5.1e-2	1.3e-1

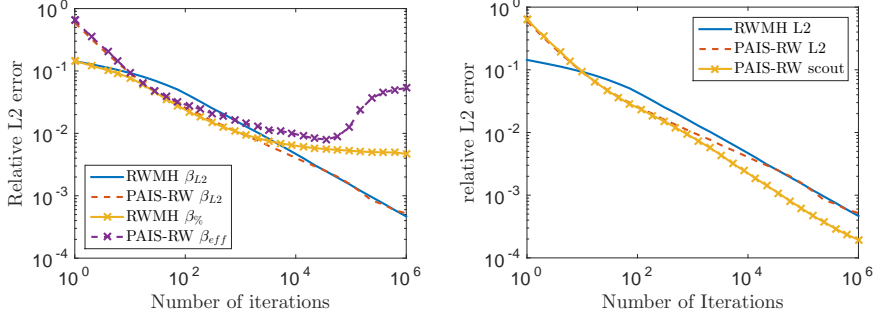
TABLE 7.2
Optimal scaling parameters for the example in Section 7.3.

Since transitions between the modes are extremely unlikely, we consider the convergence on two levels; convergence to equally proportioned modes, and convergence to smooth histograms. To get correctly proportioned modes with the RWMH algorithm it is important that the chains can transition between the modes frequently, which requires a large proposal variance. However, this leads to a lower acceptance rate, and so we sacrifice convergence locally. For this reason, the RWMH algorithm is very slow to converge for problems of this type.

We can achieve these two regimes in RWMH by tuning β using different statistics; the acceptance rate for local convergence, and the L^2 error for global convergence. Similarly in PAIS-RW we can use the effective sample size for local convergence, and the L^2 error for global convergence.

7.3.3. Convergence of RWMH vs PAIS-RW. The MH and PAIS algorithms are run with both the globally and locally optimal scaling parameters. Figure 7.5 (a) shows that the globally optimal β^* leads to normal convergence rates, ~~$\mathcal{O}(n^{-1/2})$~~ $\mathcal{O}(1/\sqrt{N})$, whereas algorithms using locally optimal scaling parameters ~~initially converge faster but~~ do not sample proportionally for the whole simulation.

Using the PAIS algorithm we can ensure that we do not lose modes as the simulation progresses. Since we have an ensemble, we can sacrifice locally optimal sampling for one or two ensemble members in favour of a larger proposal variance. These chains with larger scaling parameters act both as ‘scouts’ for new modes, and act to aid in the over-dispersal of the proposal distribution. Figure 7.5 (b) shows the results of using 49 chains with the local optimal scaling parameter, and one chain with ten times the local optimal scaling parameter. We see that modes are not forgotten and the algorithm converges with the improvement we see from PAIS-RW in the Gaussian example.



(a) Relative L^2 error using the optimal scaling parameters. (b) Convergence of locally optimised PAIS-RW with scout chain, and globally optimised PAIS-RW and RWMH.

FIG. 7.5. Convergence of the PAIS-RW and RWMH algorithms for the bimodal example discussed in Section 7.3. Set up described in Section 7.1. Resampling is performed using the ETPF.

The adaptive algorithm can just as easily be applied to the PAIS algorithm with the ‘scout’ chains described in the previous paragraph. Since we need two equally sized sub-ensembles, we will use two groups of 25 ensemble members. Each group has 24 ensemble members tuned for locally optimal convergence and one ‘scout’ chain with a scaling parameter ten times that of the rest of the group.

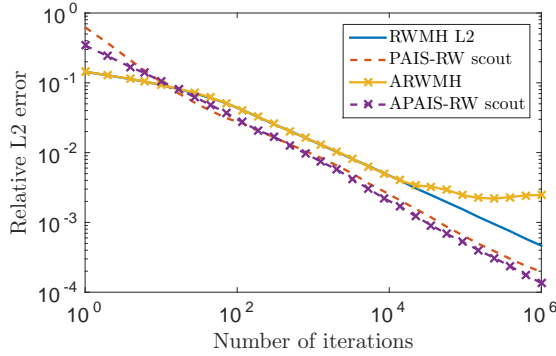


FIG. 7.6. Convergence of modified PAIS and RWMH algorithms, comparing the globally optimised nonadaptive algorithms with the locally optimised adaptive algorithms. The setup is as described in Section 7.1 (3). Resampling is performed using the ETPF.

Comparing the convergence of the adaptive algorithms against the nonadaptive algorithms in Figure 7.6 shows that the algorithms behave as expected. The adaptive RWMH algorithm tuned using the acceptance rate converges at the same rate as the L^2 optimised algorithm until the scaling parameter gets too small, and therefore switches between the modes are rare, and the relative heights of the modes are decided by the arbitrary proportion of chains which are in each mode at this point. The adaptive PAIS-RW with scout chains tuned to the effective sample size converges at about the same rate as the locally optimised nonadaptive algorithm also with scouts.

7.3.4. Calculating the speed up in convergence. The graphs in the previous section clearly show that the PAIS-RW algorithm converges faster than the RWMH

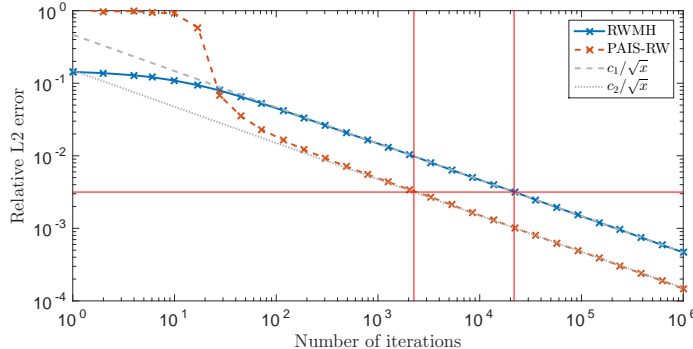


FIG. 7.7. Illustration of calculating the number of PAIS-RW iterations required to reach a tolerance of $10^{-2.8}$ as a percentage of MH iterations. The relative L^2 error graphs are from the Gaussian problem in Section 7.2.

algorithm when both are parallelised with the same ensemble size. We can calculate the number of iterations required to achieve a particular tolerance level in our solution for each algorithm and compare these to calculate a percentage saving. In Figure 7.7 we demonstrate our calculation of the savings. The constants c_1 and c_2 are found by regressing through the data with a fixed exponent of $-1/2$ excluding the initial data points where the algorithm has not finished burning in.

	Gaussian	Bimodal
RWMH	10%	12.6% (scout)
MALA	42%	40%

TABLE 7.3

Iterations for the PAIS algorithms required to achieve a desired tolerance as a percentage of the number of iterations required by the respective MH algorithms.

A summary of the percentage of iterations required using the PAIS algorithm compared with the respective Metropolis-Hastings algorithms is given in Table 7.3.

7.3.5. A useful property of the PAIS algorithm for multimodal distributions. The biggest issue for the Metropolis-Hastings algorithms when sampling from bimodal posterior distributions is that it is unlikely that the correct ratio of chains will be maintained in each of the modes, and since there is no interaction between the chains, there is no way to remedy this problem. The PAIS algorithm tackles this problem with its resampling step. The algorithm uses its dynamic [kernel-mixture proposal distribution](#) to build up an approximation of the posterior at each iteration, and then compares this to the posterior distribution via the weights function. Any large discrepancy in the approximation will result in a large or small weight being assigned to the relevant chain, meaning the chain will either pull other chains towards it or be sucked towards a chain with a larger weight. In this way, the algorithm allows chains to ‘teleport’ to regions of the posterior which are in need of more exploration. Figure 7.8 shows the trace of a simulation of the PAIS-RW algorithm with initially 1 chain in the positive mode, and 49 chains in the negative mode. It takes only a handful of iterations for the algorithm to balance out the chains into 25 chains in each mode. The chains switch modes without having to climb the energy gradient in the

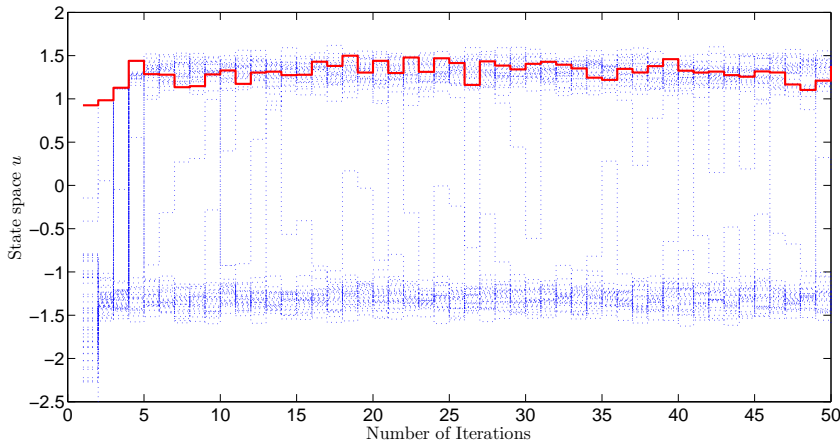


FIG. 7.8. This figure demonstrates the redistribution property of the PAIS algorithm. Initially there is one chain in the positive mode, and 49 chains in the negative mode.

middle.

7.4. A Mixture Model. The technique of mixture modelling employs well known parametric families to construct an approximating ~~nonparametric~~-distribution which may have a complex structure. Most commonly, Gaussian kernels are used since underlying properties in the data can often be assumed to follow a Gaussian distribution. An example would be if a practitioner were to measure the heights of one hundred adults, but failed to record their gender. The data could be considered as one population with two sub populations, male and female. The problem then might be to find the average height of adult males from the data. In this case, since height is often considered to follow a Gaussian distribution, it makes sense to model the population as a mixture of two univariate Gaussian distributions.

A well known problem in the Bayesian treatment of mixture modelling is that of identification, sometimes referred to as the label-switching phenomenon. The likelihood distribution for mixture models is invariant under permutations of the mixture labels. If a mixture has n means and the point (μ_1, \dots, μ_n) maximises the likelihood, then the likelihood will also be maximised by $(\mu_{\varphi(1)}, \dots, \mu_{\varphi(n)})$ for all permutations $\varphi(\cdot)$. This means that the number of modes in the posterior distribution is of order $\mathcal{O}(n!)$. As we have seen it can be hard for standard MH algorithms to obtain reliable inference for posterior distributions with a large number of modes, or even a small number of modes which are separated by a large distance.

7.4.1. Target Distribution. In particular we look at a data set where we assume that there are two subpopulations within the overall population. Since both subpopulations will be approximated by Gaussian distributions we have five parameters which we need to be estimated, two means $\{\mu_1, \mu_2\}$, two variances $\{\sigma_1^2, \sigma_2^2\}$, and the probability, p , that an individual observation belongs to the first subpopulation. We have 100 data points, D_i , which we assume to be distributed according to

$$D_i \sim p\mathcal{N}(\mu_1, \sigma_1^2) + (1 - p)\mathcal{N}(\mu_2, \sigma_2^2), \quad i = 1, \dots, 100,$$

where $p \in [0, 1]$, $\mu_1, \mu_2 \in \mathbb{R}$ and $\sigma_1^2, \sigma_2^2 \in \mathbb{R}^+$. Due to the domains of these parameters and also some prior knowledge, we assign the priors

$$p \sim \text{Beta}(1, 1), \quad \mu_{1,2} \sim \mathcal{N}(0, 4) \quad \text{and} \quad \sigma_{1,2}^2 \sim \text{Gamma}(\alpha = 2, \beta = 1).$$

If we collect these parameters in the vector $\theta = (p, \mu_1, \sigma_1^2, \mu_2, \sigma_2^2)^\top$, the resulting posterior distribution is

$$\pi(\theta|D) \propto \prod_{i=1}^{100} (p\mathcal{N}(D_i; \mu_1, \sigma_1^2) + (1-p)\mathcal{N}(D_i; \mu_2, \sigma_2^2)) \prod_{i=1}^5 \pi_0^i(\theta_i), \quad (7.2)$$

where $\pi_0^i(\cdot)$ is the prior density function corresponding to θ_i .

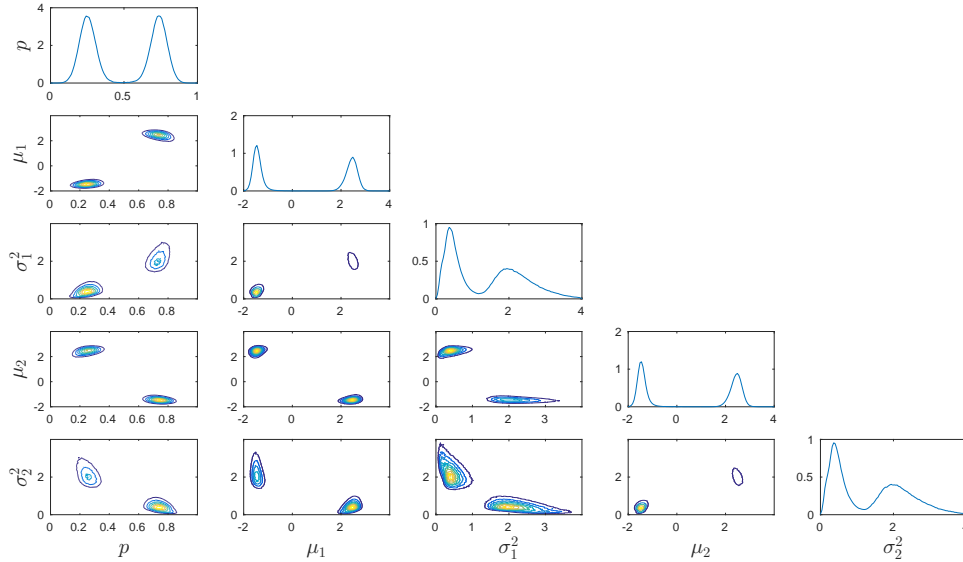


FIG. 7.9. The posterior distribution of θ as given in Equation (7.2) as found from 10 million samples from the PAIS algorithm. The main diagonal contains the marginal distributions of each θ_i , and the lower triangular contours represents the correlation between pairs of parameters.

Figure 7.9 presents a visualisation of the posterior distribution for this problem, created from ~~10million~~ 10 million samples produced by the PAIS-RW algorithm.

7.4.2. Implementation. We have no analytic form for the posterior distribution for this model, and MH cannot reliably sample from this type of bimodal distribution. We do however know that the probability density should be evenly divided between the two modes. This is due to the symmetric prior for p , and the same priors being assigned to μ_1 and μ_2 , and also to σ_1^2 and σ_2^2 . To decide which mode a sample belongs to we define a plane which bisects the posterior so that each point on this plane lies exactly halfway between the two true solutions to the inverse problem i.e. the value of θ used to generate the data, and also the θ obtained by a relabelling of the parameters. Now that we can assign a sample to a particular mode, we can calculate the density in each mode by summing the weights associated to all samples

in that mode,

$$\bar{w}_k = \sum_{i=1}^N w_i I_{X_i \in \text{Mode } k}, \quad k = 1, 2,$$

and the relative error in the amount of density in each mode is then

$$w_{\text{error}} = 2 \left| \frac{\bar{w}_1}{\bar{w}_1 + \bar{w}_2} - \frac{1}{2} \right|. \quad (7.3)$$

Since the probability p is constrained to lie in the interval $[0, 1]$, and the variances must be positive, it can be wasteful to use Gaussian proposal distributions. For this example we make proposals by matching the proposal distribution to the prior, scaling the variance and centring the proposal at the previous value. So

$$q_p \sim \text{Beta}(\delta^{-2}p, \delta^{-2}(1-p)), \quad q_{\mu_{1,2}} \sim \mathcal{N}(\mu_{1,2}, 4\delta^2) \quad \text{and} \quad q_{\sigma_{1,2}^2} \sim \text{Gamma}(\alpha^*, \beta^*),$$

where $\alpha^* = \sigma_{1,2}^2 \beta^*$, $\beta^* = \sigma_{1,2}^2 / 2\delta^2$ and δ is a scaling parameter to be tuned. This means that our proposal distributions will not be a mixture of multivariate Gaussians, but independent mixtures of univariate Beta, Gamma and Gaussian distributions.

In the numerics which follow we have increased the ensemble size from $M = 50$ to $M = 500$ to compensate for the increase in dimension. We also use the AMR algorithm for the resampling step because of the reduced computational cost. The method otherwise remains the same as in previous examples. We perform test runs to find the optimal scaling parameters considering convergence to modes with equal density. We then calculate the convergence rates of the algorithms by producing 10 million samples from the posterior with each algorithm, and repeat the simulation 32 times.

Algorithm	MH	PAIS
δ^*	1.3e-1	2.3e-1

TABLE 7.4

Optimal values of the scaling parameter. The MH algorithm is optimised using the acceptance rate, and the PAIS algorithm is optimised using the effective sample size.

7.4.3. Convergence of MH vs PAIS. The optimal scaling parameters for the MH and PAIS algorithms with the proposal distributions described in Section 7.4.2 are given in Table 7.4.

Convergence of the relative error for the two algorithms is displayed in Figure 7.10. PAIS converges at the expected $\mathcal{O}(1/\sqrt{N})$ rate, whereas the MH algorithm converges to locally smooth histograms but with the wrong proportion of samples in each mode. The relatively low value of the error for the MH example is due to the priors covering the sample space evenly, however since transitions are near impossible with a small value of the scaling parameter, this error will take a very long time to reduce. This problem was discussed in Section 7.3.

In this example, we have only considered a relatively low dimensional mixture model problem, with only two distributions in the mixture. With more elements in the mixture, and/or an undefined number in the mixture, the dimension of this problem will quickly increase. Importance sampling schemes such as this suffer from the curse of dimensionality, limiting the size of the dimension of target density that it can efficiently sample from. However, this example does demonstrate the remarkable convergence of the PAIS for multimodal target distributions.

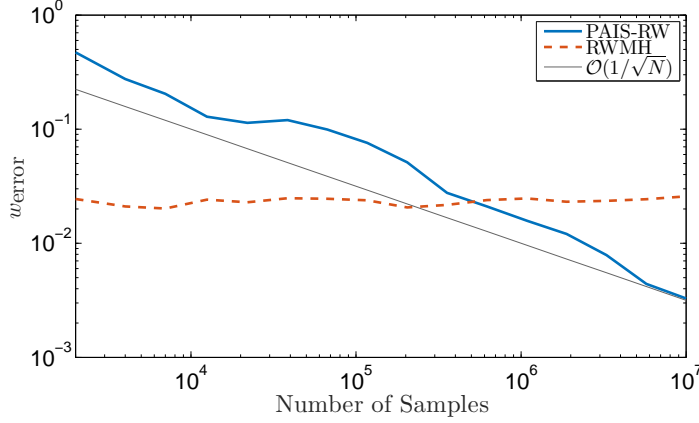


FIG. 7.10. Convergence of the PAIS algorithm for the mixture model described in Section 7.4, convergence judged using the criterion in Equation (7.3). Implementation described in Sections 7.4.2 and 7.4.3. Resampling is performed using the AMR scheme.

7.5. Data assimilation with Lorenz ‘63 trajectories. The sampling algorithm we have introduced in this paper is designed for problems where the likelihood density function is expensive to calculate, and the computational overhead involved in the calculation of the importance weights is dwarfed. This situation commonly occurs in inverse problems where the model being investigated involves the solution of a differential equation. In this example we attempt to recover the initial position of a particle with motion governed by Lorenz’s 1963 atmospheric convection equations [29]:

$$\begin{aligned}\frac{dx}{dt} &= \sigma(y - x), \\ \frac{dy}{dt} &= x(\rho - z) - y, \\ \frac{dz}{dt} &= xy - \beta z.\end{aligned}$$

When the parameters are chosen to be $\rho = 28$, $\sigma = 10$ and $\beta = 8/3$ this system has chaotic solutions. If we are interesting in calculating the initial position by observing the location of the particle at certain points in time, it can quickly become intractable when there is any noise in the observations.

To demonstrate how small errors in the initial condition of a Lorenz solution can affect the trajectory of a particle, Figure 7.11 shows the path taken by two particles with very similar initial conditions,

$$\mathbf{x}_0^1 = \begin{pmatrix} -0.587 \\ -0.563 \\ 16.870 \end{pmatrix} \quad \text{and} \quad \mathbf{x}_0^2 = \begin{pmatrix} -0.590 \\ -0.563 \\ 16.870 \end{pmatrix}.$$

This small difference in the x -dimension of the initial condition leads to the trajectories of the two particles decoupling near $t = 5$. This chaotic behaviour means that if we

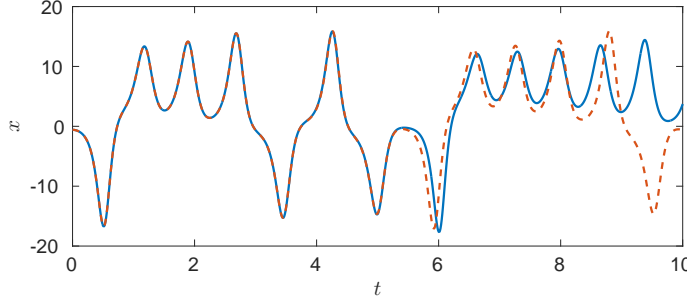


FIG. 7.11. *Trajectory of the x component of a particles position varying in time when motion is governed by the Lorenz '63 equations. The two trajectories are calculated using an initial position of $[-0.587, -0.563, 16.870]$ and $[-0.590, -0.563, 16.870]$.*

take only a few observations over a long time period we will find that many trajectories which may be vastly different achieve similar values of the likelihood function and so the posterior distribution for the initial condition can become very complex.

7.5.1. Target Distribution. For this example, we observe noisily the position of a particle at ten equally spaced points in the time interval $t \in (0, 1]$. The Lorenz equations are solved with the chaotic parameters given above and the initial condition \mathbf{x}_0 . A time step of $h = 1 \times 10^{-3}$ was used to evolve the equations numerically with the explicit Euler method. The noise added to each observation was taken from the distribution $\mathcal{N}(0, 0.1^2 \mathbf{I})$. Priors for the initial condition coordinates were given by

$$x_0 \sim \mathcal{N}(-0.5, 0.03^2), \quad y_0 \sim \mathcal{N}(-0.5, 0.03^2) \quad \text{and} \quad z_0 \sim \mathcal{N}(15, 0.03^2).$$

The posterior density function takes the form

$$\pi(\mathbf{x}_0 | D) \propto \exp \left\{ -\frac{1}{2} \|\mathcal{G}(\mathbf{x}_0) - D\|_{\Sigma}^2 \right\} \mathcal{N}(\mathbf{x}_0; m, 0.03^2 \mathbf{I}),$$

where \mathcal{G} calculates the solution to the Lorenz equations, using the Euler method with a time step of $h = 1 \times 10^{-3}$ starting at the initial condition \mathbf{x}_0 , and returns the position at $t = 0.1, \dots, 1$. The mean of the prior distribution $m = [-0.5, -0.5, 15]$. Figure 7.12 shows the marginal distributions of the posterior distribution for the initial condition \mathbf{x}_0 as found using the MH algorithm. While the z -dimension is largely uncorrelated with the other two dimensions, the correlation between x and y is -0.98, which makes it a challenging posterior distribution to explore without a transformation of the parameter space.

7.5.2. Implementation. As in the previous example, we have no analytic form for the normalisation constant of the posterior and so we will measure convergence of the algorithms by calculating convergence to the posterior mean, where the truth is calculated using a very long chain produced by the MH algorithm. Both the MH and PAIS algorithms use an ensemble size of $M = 500$. The PAIS algorithm uses the AMR resampling algorithm. The optimal scaling parameters are calculated for both algorithms using test runs which are not included in the convergence cost calculations. Convergence graphs are produced using 32 repeats of each algorithm with each simulation producing one million samples.

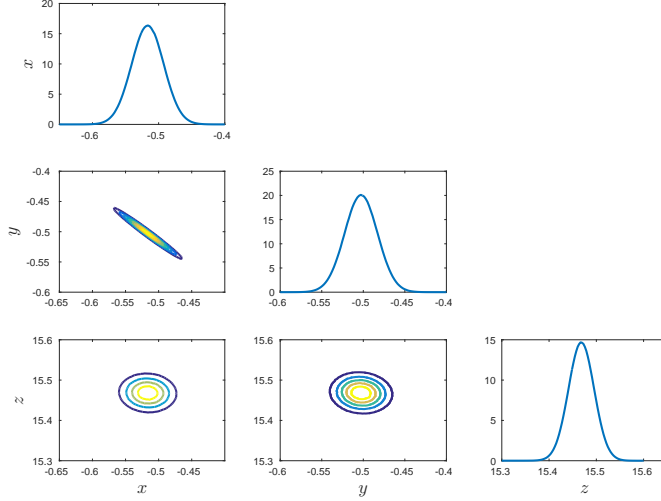


FIG. 7.12. *Posterior distribution for the initial condition of the particle with motion governed by the Lorenz equations as described in Section 7.5.1.*

7.5.3. Convergence of MH vs PAIS. The optimal scaling parameters for the MH and PAIS algorithms are given in Table 7.5.

Algorithm	MH	PAIS
δ^*	1.0e-1	1.3e-2

TABLE 7.5

Optimal values of the scaling parameter. The MH algorithm is optimised using the acceptance rate, and the PAIS algorithm is optimised using the effective sample size.

As is clear from Figure 7.13 the PAIS-RW has a vastly superior burn-in time and achieves improved accuracy on the RWMH algorithm of three orders of magnitude after one million samples.

To produce one sample from the posterior distribution for this example using the PAIS algorithm takes on average 5.5×10^{-4} seconds, while a sample can be produced in 1.7×10^{-5} seconds using the MH algorithm. However; the quality of the sample produced by PAIS is much higher. To reach a tolerance of 5×10^{-5} , it takes the PAIS algorithm roughly 16.5 seconds and the MH algorithm roughly 1020 seconds. Although we do not claim that either codes are optimally efficient, this gives an indication of the saving that can be achieved through this type of methodology.

8. Discussion and Conclusions. We have explored the application of parallelised MCMC algorithms in low dimensional inverse problems. We have demonstrated numerically that these algorithms converge faster than the analogous naively parallelised Metropolis-Hastings algorithms. Further experimentation with the Metropolis Adjusted Langevin Algorithm (MALA), preconditioned Crank-Nicolson (pCN), preconditioned Crank-Nicolson Langevin (pCNL) and Hamiltonian Monte Carlo (HMC) proposals has yielded similar results [44]. Importantly, we have compared the efficiency of our parallel scheme with a naive

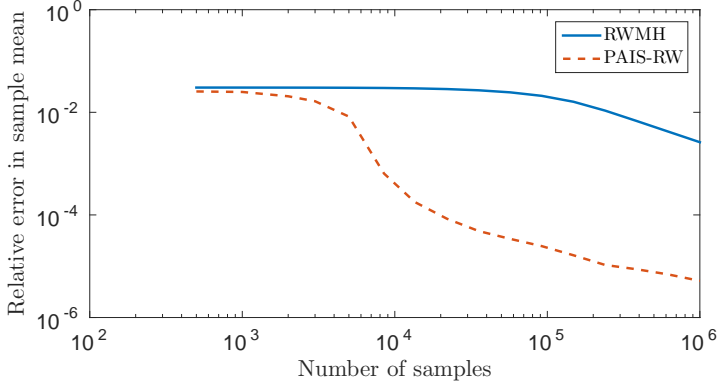


FIG. 7.13. *Convergence of the sample mean to the posterior mean for the MH and PAIS algorithms. Posterior distribution as described in Section 7.5.1.*

parallelisation of serial methods. Thus our increase in efficiency is over and above an M -fold increase, where M is the number of ensemble members used. Our approach demonstrates a better-than-linear speed-up with the number of ensemble members used.

The PAIS has a number of favourable features, for example the algorithm’s ability to redistribute, through the resampling regime, the ensemble members to regions which require more exploration. This allows the method to be used to sample from complex multimodal ~~distribution~~distributions.

Another strength of the PAIS is that it can also be used with any MCMC proposal. There are a growing number of increasingly sophisticated MCMC algorithms (non-reversible Langevin/HMC proposals, Riemann manifold MCMC etc) which could be incorporated into this framework, leading to even more efficient algorithms, and this is another opportunity for future work.

One limitation of the PAIS approach as described above is that a direct solver of the ETPF problem (such as FastEMD [36]) has computational cost $\mathcal{O}(M^3 \log M)$, where M is the number of particles in the ensemble. As such, we introduced a more approximate resampler the approximate multinomial resampler, which allows us to push the approach to the limit with much larger ensemble sizes. The PAIS framework is very flexible in terms of being able to use any combination of proposal distributions and resampling algorithms that one wishes.

We have demonstrated that the framework that we have considered, with the use of state-of-the-art optimal transport-based resampling, can reduce the number of likelihood evaluations required to characterise complex posterior distributions in low dimensions to a given degree. We have also introduced a greedy approximation to this resampler, which drastically reduces the cost, at the loss of some accuracy, which can be mitigated with the use of larger ensembles. We have detailed how scaling parameters in the MCMC proposals that are used within the mixture distribution can be quickly and efficiently tuned in an automated way. Lastly we have demonstrated that when the likelihood is expensive, for instance because it involves the numerical approximation of the solution of a differential equation, the savings made in terms of the number iterations required in order reach a particular degree of accuracy outweighs the additional overhead costs that are incurred over-and-above standard MCMC.

However, tuning the variances and covariances of the mixture proposal components globally is likely to be of limited use for multimodal problems where the modes have very different covariances, or where there is curved ridges in the density. There are a variety of possible solutions to this problem that would be worthy of future consideration, not least the potential to use transport maps [18, 35] to map Gaussian mixtures to highly complex non-Gaussian approximations of the posterior. Such a map could encode local covariance information, and lead to accelerated and stable PAIS sampling. This is an interesting area for future study.

Appendix A. Glossary of acronyms.

Acronym	Full name
MCMC	Markov chain Monte Carlo
RWMH	Random walk Metropolis-Hastings
MALA	Metropolis-adjusted Langevin algorithm
AIS	Adaptive importance sampling
PMC	Population Monte Carlo
PAIS	Parallel adaptive importance sampling
ETPF	Ensemble transport particle filter
AMR	Approximate multinomial resampler
PAIS-X	PAIS with X kernels
AP AIS-X	Adaptive PAIS with X kernels

REFERENCES

- [1] C. ANDRIEU, É. MOULINES, ET AL., *On the ergodicity properties of some adaptive MCMC algorithms*, The Annals of Applied Probability, 16 (2006), pp. 1462–1505.
- [2] E. ANGELINO, E. KOHLER, A. WATERLAND, M. SELTZER, AND R. ADAMS, *Accelerating MCMC via parallel predictive prefetching*, arXiv preprint arXiv:1403.7265, (2014).
- [3] J. BIERKENS, *Non-reversible Metropolis-Hastings*, Statistics and Computing, (2015), pp. 1–16.
- [4] M. BINK, M. BOER, C. TER BRAAK, J. JANSEN, R. VOORRIPS, AND W. VAN DE WEG, *Bayesian analysis of complex traits in pedigreed plant populations*, Euphytica, 161 (2008), pp. 85–96.
- [5] M. BUGALLO, L. MARTINO, AND J. CORANDER, *Adaptive importance sampling in signal processing*, Digital Signal Processing, 47 (2015), pp. 36–49.
- [6] T. BUI-THANH AND M. GIROLAMI, *Solving large-scale PDE-constrained Bayesian inverse problems with Riemann manifold Hamiltonian Monte Carlo*, Inverse Problems, 30 (2014), p. 114014.
- [7] B. CALDERHEAD, *A general construction for parallelizing Metropolis-Hastings algorithms*, Proceedings of the National Academy of Sciences, 111 (2014), pp. 17408–17413.
- [8] OLIVIER CAPPÉ, RANDAL DOUC, ARNAUD GUILLIN, JEAN-MICHEL MARIN, AND CHRISTIAN P ROBERT, *Adaptive importance sampling in general mixture classes*, Statistics and Computing, 18 (2008), pp. 447–459.
- [9] O. CAPPÉ, A. GUILLIN, J. MARIN, AND C. ROBERT, *Population Monte Carlo for ion channel restoration*.
- [10] O. CAPPÉ, A. GUILLIN, J. MARIN, AND C. ROBERT, *Population monte carlo*, Journal of Computational and Graphical Statistics, (2012).
- [11] G. CELEUX, J. MARIN, AND C. ROBERT, *Iterated importance sampling in missing data problems*, Computational statistics & data analysis, 50 (2006), pp. 3386–3404.
- [12] J. CORNUET, J. MARIN, A. MIRA, AND C. ROBERT, *Adaptive multiple importance sampling*, Scandinavian Journal of Statistics, 39 (2012), pp. 798–812.
- [13] C. COTTER AND S. REICH, *Ensemble filter techniques for intermittent data assimilation-a survey*, in Large Scale Inverse Problems. Computational Methods and Applications in the Earth Sciences, Walter de Gruyter, Berlin, 2012, pp. 91–134.
- [14] S. COTTER, M. DASHTI, J. ROBINSON, AND A. STUART, *Bayesian inverse problems for functions and applications to fluid mechanics*, Inverse Problems, 25 (2009), p. 115008.

- [15] S. COTTER, G. ROBERTS, A. STUART, AND D. WHITE, *MCMC methods for functions: modifying old algorithms to make them faster*, Statistical Science, 28 (2013), pp. 424–446.
- [16] R. DOUC, A. GUILLIN, J. MARIN, C. ROBERT, ET AL., *Convergence of adaptive mixtures of importance sampling schemes*, The Annals of Statistics, 35 (2007), pp. 420–448.
- [17] RANDAL DOUC, ARNAUD GUILLIN, J-M MARIN, AND CHRISTIAN P ROBERT, *Minimum variance importance sampling via population monte carlo*, ESAIM: Probability and Statistics, 11 (2007), pp. 427–447.
- [18] T. EL MOSELHY AND Y. MARZOUK, *Bayesian inference with optimal maps*, Journal of Computational Physics, 231 (2012), pp. 7815–7850.
- [19] VÍCTOR ELVIRA, LUCA MARTINO, DAVID LUENGO, AND MÓNICA F BUGALLO, *Improving population monte carlo: Alternative weighting and resampling schemes*, Signal Processing, 131 (2017), pp. 77–91.
- [20] G. EVENSEN, *Sequential data assimilation with a nonlinear quasi-geostrophic model using Monte Carlo methods to forecast error statistics*, Journal of Geophysical Research: Oceans (1978–2012), 99 (1994), pp. 10143–10162.
- [21] M. GIROLAMI AND B. CALDERHEAD, *Riemann manifold Langevin and Hamiltonian Monte Carlo methods*, Journal of the Royal Statistical Society: Series B (Statistical Methodology), 73 (2011), pp. 123–214.
- [22] H. HAARIO, E. SAKSMAN, AND J. TAMMINEN, *Componentwise adaptation for high dimensional MCMC*, Computational Statistics, 20 (2005), pp. 265–273.
- [23] W. HASTINGS, *Monte Carlo sampling methods using Markov chains and their applications*, Biometrika, 57 (1970), pp. 97–109.
- [24] T. HOUSE, A. FORD, S. LAN, S. BILSON, E. BUCKINGHAM-JEFFERY, AND M. GIROLAMI, *Bayesian uncertainty quantification for transmissibility of influenza, norovirus and ebola using information geometry*, Journal of The Royal Society Interface, 13 (2016), p. 20160279.
- [25] M. ISARD AND A. BLAKE, *Condensationconditional density propagation for visual tracking*, International journal of computer vision, 29 (1998), pp. 5–28.
- [26] C. JI AND S. C. SCHMIDLER, *Adaptive Markov Chain Monte Carlo for Bayesian Variable Selection*, Journal of Computational and Graphical Statistics, 22 (2013), pp. 708–728.
- [27] J. LIU, *Monte Carlo strategies in scientific computing*, Springer Science & Business Media, 2008.
- [28] J. LIU, F. LIANG, AND W. WONG, *The multiple-try method and local optimization in Metropolis sampling*, Journal of the American Statistical Association, 95 (2000), pp. 121–134.
- [29] EDWARD N LORENZ, *Deterministic nonperiodic flow*, Journal of the atmospheric sciences, 20 (1963), pp. 130–141.
- [30] J.-M. MARIN, K. MENGENSEN, AND C. ROBERT, *Bayesian modelling and inference on mixtures of distributions*, Handbook of statistics, 25 (2005), pp. 459–507.
- [31] L. MARTINO, V. ELVIRA, D. LUENGO, AND J. CORANDER, *An adaptive population importance sampler: Learning from uncertainty*, IEEE Transactions on Signal Processing, 63 (2015), pp. 4422–4437.
- [32] LUCA MARTINO, VÍCTOR ELVIRA, DAVID LUENGO, AND JUKKA CORANDER, *Layered adaptive importance sampling*, Statistics and Computing, 27 (2017), pp. 599–623.
- [33] G. MOORE ET AL., *Cramming more components onto integrated circuits*, Proceedings of the IEEE, 86 (1998), pp. 82–85.
- [34] R. NEAL, *MCMC using ensembles of states for problems with fast and slow variables such as Gaussian process regression*, arXiv preprint arXiv:1101.0387, (2011).
- [35] MATTHEW PARNO AND YOUSSEF MARZOUK, *Transport map accelerated markov chain monte carlo*, arXiv preprint arXiv:1412.5492, (2014).
- [36] O. PELE AND M. WERMAN, *Fast and robust Earth mover’s distances*, in Computer Vision, 2009 IEEE 12th International Conference on, IEEE, 2009, pp. 460–467.
- [37] S. REICH, *A nonparametric ensemble transform method for Bayesian inference*, SIAM Journal on Scientific Computing, 35 (2013), pp. A2013–A2024.
- [38] C. ROBERT AND G. CASELLA, *Monte Carlo statistical methods*, Springer Science & Business Media, 2013.
- [39] G. ROBERTS AND J. ROSENTHAL, *Optimal scaling of discrete approximations to Langevin diffusions*, Journal of the Royal Statistical Society: Series B (Statistical Methodology), 60 (1998), pp. 255–268.
- [40] ———, *Optimal scaling for various Metropolis-Hastings algorithms*, Statistical science, 16 (2001), pp. 351–367.
- [41] ———, *Coupling and ergodicity of adaptive Markov chain Monte Carlo algorithms*, Journal of Applied Probability, 44 (2007), pp. 458–475.
- [42] ———, *Examples of adaptive MCMC*, Journal of Computational and Graphical Statistics, 18

- (2009), pp. 349–367.
- [43] J. ROSENTHAL, *Optimal proposal distributions and adaptive MCMC*, Handbook of Markov Chain Monte Carlo, (2011), pp. 93–112.
 - [44] P. RUSSELL, *PhD Thesis*, PhD thesis, School of Mathematics, University of Manchester, 2017.
 - [45] J. SEXTON AND D. WEINGARTEN, *Hamiltonian evolution for the hybrid Monte Carlo algorithm*, Nuclear Physics B, 380 (1992), pp. 665–677.
 - [46] J. SIRÉN, P. MARTTINEN, AND J. CORANDER, *Reconstructing population histories from single nucleotide polymorphism data*, Molecular biology and evolution, 28 (2011), pp. 673–683.
 - [47] A. SOLONEN, P. OLLINAHO, M. LAINE, H. HAARIO, J. TAMMINEN, H. JÄRVINEN, ET AL., *Efficient MCMC for climate model parameter estimation: parallel adaptive chains and early rejection*, Bayesian Analysis, 7 (2012), pp. 715–736.
 - [48] A. STUART, *Inverse problems: a Bayesian perspective*, Acta Numerica, 19 (2010), pp. 451–559.