

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 MCMC sampling algorithms. At each iteration, an importance sampling proposal distribution is formed using the current states of all of the chains within an ensemble. Once weighted samples have been produced from this 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 using the same number of ensemble members, 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. In particular, we demonstrate this methodology's superior sampling for multimodal problems, such as those arising from inference for mixture models.

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

1. Introduction. Having first been developed in the early 1970s [21], Markov chain Monte Carlo (MCMC) methods are a powerful family of tools that 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. MCMC methods were first developed in the 70s [21], and with the development of faster, more powerful computers, have become ever more important in a whole range of fields in statistics, science and engineering. In particular, when considering Bayesian inverse problems, each MCMC step may involve the numerical solution of one or more PDE. As many samples are usually required before Monte Carlo error is reduced to acceptable levels, the application of MCMC methods to these types of problem remain frustratingly out of our grasp.

Many advances have been made in the field of MCMC to design ever more complex methods that propose moves more intelligently. 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 rapidly converging approximations. Function space versions of standard methods such as the random walk Metropolis-Hastings (RWMH) algorithm or the Metropolis-adjusted Langevin algorithm (MALA), whose convergence rates

*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

are independent of dimension have been developed [14]. The hybrid (or Hamiltonian) more thorough exploration of the posterior with fewer iterations. For example, the Hamiltonian or Hybrid Monte Carlo (HMC) method uses Hamiltonian dynamics algorithm uses gradient information and symplectic integrators in order to propose and accept moves to states which are a long way away from the current position [41], and function space analogues of this have also been proposed [2] make very large moves in state with relatively high acceptance probability [41]. 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 [18]. This methodology has been successfully applied to MALA-type proposals and methods which exploit even higher order gradient information [6]. These methods allow us to explore the posterior distribution more fully with fewer iterations.

Simultaneously, great strides are continually being made in the development of computing hardware. Since the clock speed of an individual processor is no longer following Moore's law, which predicted that the number of transistors that can fit onto a single microchip will double every two years, has been largely followed since the early 70s [31]. In recent times, it has become necessary to use parallel architectures in order for this trend to continue. The efficient exploitation of these architectures is the key to solving many of the computational challenges that we currently face.

[31], improvements in computational power are largely coming from the parallelisation of multiple cores. As such, the development of efficient parallel MCMC algorithms is an important area for research. Since MCMC methods can be naively parallelised by simply running many independent chains in parallel, the focus needs to be on the development of methods which gain some added benefit through parallelism 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 [32]. 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.

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 [14]. Other ideas such as chain adaptation [20] and early rejection of samples can also aid reduction of the computational workload [43].

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 [22]. 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 [29]. 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, 8, 10, 23]. 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 [30, 35].

Ensemble importance sampling schemes also exist, e.g. population Monte Carlo (PMC) [9]. 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. Extensions, such as adaptive multiple importance sampling algorithm (AMIS) [11], and adaptive population importance sampling (APIS) [30] have enabled these methods to be applied to various applications, including population genetics [42].

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 ~~proposal distribution~~ from the current ~~state of all of the chains in the ensemble~~. The ~~idea is to consider the current state on each of a set of parallel chains as an ensemble, and to resample using a transformation method makes use of a resampler~~ based on optimal transport ~~. Samplers based on optimal transport have also been considered in [15], which has been used in the context of particle filters [34].~~ 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.

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 parameter quantity x is ~~observed through mapped on to observable space by~~ the observation operator $\mathcal{G} : X \rightarrow \mathbb{R}^d$. ~~Since observations are never perfect, we assume that these measurements~~ We assume that the observations, D , are subject to Gaussian noise, ~~so that~~

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

For example, if x ~~are the rates of reactions in a chemical system,~~ \mathcal{G} might return the quantities of each chemical species is an initial condition for the wave equation, \mathcal{G} might calculate the height of a wave at a particular ~~time, or some summary of this information point and time.~~

These modelling assumptions allow us to construct the likelihood of observing the data D given the parameter 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 ~~$\|y_1 - y_2\|_\Sigma$ is the Mahalanobis distance between y_1 and $y_2 \in \mathbb{R}^d$~~ $\|y_1 - y_2\|_\Sigma = (y_1 - y_2)^\top \Sigma^{-1} (y_1 - y_2)$ for $y_1, y_2 \in \mathbb{R}^d$

As discussed in [13, 44], in order for this inverse problem to be well-posed in the Bayesian sense, we require the posterior distribution, 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:

$$\frac{d\mu_Y}{d\mu_0} \propto \exp(-\Phi(x^*)). \quad (2.3)$$

2.2. Particle filters and resamplers. ~~In several applications, data must be assimilated in an “online” fashion, with up to date observations of the studied system being made available on a regular basis. In these contexts, such as in weather forecasting or oceanography, data is incorporated using a filtering methodology. One popular filtering method is the particle filter, 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 first of which was dubbed the Bootstrap filter [19]. In this method, a set of weighted particles is used to represent the posterior distribution. The positions of the particles are updated using the model dynamics. Then, when more observations are made available, the relative weights of true state of a system when given only noisy observations of the particles are updated to take account of this data, using Bayes’ formula. Other filtering methods, such as the Kalman filter [25] and ensemble Kalman filter [16], have also been developed which are often used within 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 — One advantage of the particle filter is that there are convergence results for this method as the number of particles is increased. One downside is that the required ensemble size increases quickly with dimension, making it difficult to use in high-dimensional problems. Another downside is that the effective sample size decreases at each iteration, resulting in degeneration of the approximation of the posterior. One way to tackle this is to employ a resampling scheme. The aim of a successful resample is to take the unevenly weighted ensemble and return a new ensemble of particles with even weights which is highly correlated to the original samples.

The has contributed several efficient particle filters, including the ensemble Kalman filter (EnKF) [16] and the ensemble transform particle filter (ETPF) proposed by Reich [34] makes use of optimal transportation as described in [45, 46]. The transform takes a sample of weighted particles $\{y_i\}_{i=1}^M$ from μ_Y and converts it into a sample of evenly weighted particles $\{x_i\}_{i=1}^M$ from μ_X , by means of defining a coupling T^* between [34].

The ETPF defines a coupling T between two random variables Y and X . Given that a trivial coupling T^I always exists in the space of transference plans, $\Pi(\mu_X, \mu_Y)$, we can find a coupling, allowing us to use the induced map as a resampler. An optimal coupling T^* is one which maximises the correlation between X and Y [12]. 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 ensures that the new sample is as much like the original sample as possible with the additional property that the sample is evenly weighted preserves the statistics of X in the new sample.

A Monte Carlo algorithm can be implemented to resample from a weighted ensemble. We create a weighted sample, then solve the optimal transport problem which produces the coupling described above, we can draw a new sample from the evenly weighted distribution. Reich suggests using the mean of the evenly weighted distribution to produce a consistent estimator. Analysis of this method shows that as the In this work we use 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, $\hat{\eta} \approx \mu_Y$.

$$\hat{\mu}_Y = \sum_{i=1}^M w_i \delta_{y_i}(\cdot) \xrightarrow{\text{ETPF}} \sum_{i=1}^M n_i \delta_{x_i}(\cdot) = \hat{\eta}. \quad (2.4)$$

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 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, the statistics of the evenly weighted sample approach those of the posterior distribution 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 over a set of processors in an ensemble 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 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 ~~,- each coming from a MCMC chain,-~~ 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 ~~distributions for each of the members of the ensemble.~~ kernels. This corresponds to replacing the Dirac distribution in the right-hand side of Equation 2.4 with a chosen distribution e.g. Gaussian. Once the proposal is constructed, we can sample a new ~~set of states ensemble~~ from the proposal distribution, and each is assigned a weight given by the ratio of the target density and the proposal mixture ~~distribution~~ 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 distribution for the next iteration of the algorithm.

Initially we will use the ETPF algorithm [34], although we will suggest an alternative strategy in Section 5. The ~~output of the~~ resampling algorithm gives us a set of evenly weighted samples ~~that we believe which~~ represents the target distribution well, ~~and from this point we can which we can use to~~ iterate the process ~~once~~ again. The algorithm is summarised in ~~Table ??~~ 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 $\mathbf{X} = [x_1, x_2, \dots, x_M]^T$ ~~$\mathbf{X} = [x_1, x_2, \dots, x_M]^T$~~ . 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.

~~$$\mathbf{X}^{(0)} = \mathbf{X}_0 = [x_1^{(0)}, x_2^{(0)}, \dots, x_M^{(0)}]^T \quad \mathbf{Y}^{(i)} = [y_1^{(i)}, y_2^{(i)}, \dots, y_M^{(i)}]^T, \quad y_j^{(i)} \sim \nu(\cdot; x_j^{(i-1)}) \quad \chi(y; \mathbf{X}^{(i)}) = \frac{1}{M} \sum_{j=1}^M \nu(y; x_j^{(i)}) \quad \mathbf{W}^{(i)} = [w_1^{(i)}, \dots, w_M^{(i)}]^T$$~~

Algorithm 1: ~~A pseudo-code representation of the Parallel Adaptive Importance Sampler~~ The parallel adaptive importance sampler (PAIS).

```

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

4   Resample:  $(\mathbf{W}^{(i)}, \mathbf{Y}^{(i)}) \rightarrow (\frac{1}{M} \mathbf{1}, \mathbf{X}^{(i)})$ 
5 Output  $(\mathbf{W}, \mathbf{Y})$ .
```

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

The key is to choose a suitable transition kernel ν such that if $X^{(i)}$ is a ~~good~~ 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 ~~RWMH~~ 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 ~~MCMC-MH~~ algorithm and lead to ~~convergence~~ 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 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 ~~simple examples~~ posteriors that we have considered [40].

Care needs to be taken when choosing the proposal distribution to ensure that the proposals are absolutely continuous with respect to the posterior distribution. ~~In Section ?? we will consider an inverse problem with Gamma priors. Since the posterior distribution in this context is heavy-tailed, if we use a lighter tailed distribution for ν , such as a Gaussian kernel, then importance weights~~ If this is not the case, importance weights grow exponentially for samples in the tails will be unbounded, which will hamper convergence of the algorithm. This will drastically decrease the efficiency of sampling.

4. Automated tuning of ~~Algorithm Parameters~~ 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 . ~~If the proposal distribution is too over-dispersed, then the algorithm will often propose states in the tails of the distribution, resulting in larger variance of the weights, and therefore slower convergence to the posterior distribution. Similarly, if the proposal distribution is under-dispersed, the proposals will be highly correlated with the previous states, and the algorithm will take a long time to fully explore the parameter space, and worse, will lead regularly to states proposed in the tails of the proposal distribution with very large weights. It is therefore necessary to find a proposal distribution which is slightly over-dispersed to ensure the entire posterior is explored [17], but is as close to the posterior as possible.~~

Most ~~commonly used~~ MCMC proposals have parametric dependence which allows the user to control their variance. For example, in the ~~RWMH~~ RW proposal $y = x + \beta\eta$, the parameter β ~~controls how correlated the proposal state y is to the current state x~~ is the standard deviation of the proposal distribution. Therefore the proposal distributions can be tuned such that they are slightly over-dispersed, ~~as described above~~. 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 [37, 38]. 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 ~~the burn-in~~ early iterations has some added benefits over and above finding an optimal parameter regime for the algorithm. If the initial value of the proposal ~~variances~~ variance is chosen to be very large, then ~~proposed moves will be made far and wide, expediting 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. ~~The fact that we have~~ Using an ensemble of chains allows ~~us to assess quickly and effectively what the optimal variance of the proposal distributions should be.~~ quick and effective assessment of the value of the optimal scaling parameter.

~~The alternative to using adaptive procedures to tune the scaling parameters is to~~

perform exploratory simulations to find the optimal regimes by trial and error. This can be very costly and the optimal parameters cannot realistically be found to more than a couple of significant figures. It is therefore important that MCMC algorithms provide a feasible adaptive strategy.

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 [36]. Unlike Metropolis-Hastings algorithms, the PAIS algorithm does not accept or reject proposed values, so we need another method of measuring the optimality of β . 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. MCMC algorithms can be assessed by comparing their approximation of the posterior to 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 analytic distribution, in cases where the posterior distribution can be computed using alternative methods. To assess this, a distance metric on distributions must be chosen. Examples are the relative error between the sample moments and the posterior target's moments, or the relative L^2 error between the true sample histogram and the target density, $\pi(x|D)$, and the constructed histogram. The relative error in the n -th moment, \hat{m}_n , is given by:

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

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

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

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

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 knowledge of the analytic solution, and the algorithm must be run for a long time to build up a sufficiently large sample are expensive to approximate. However they can be used to assess the ability of other indicators to find the optimal proposal variances scaling parameters in a controlled setting. The following statistics can be used specifically for importance sampling algorithms.

4.1.2. The variance of the weights. Importance samplers assign a weight to each sample they produce based on a ratio of the posterior to the proposal at that point. Importance samplers are most efficient when the target is proportional to the proposal distribution. In this case the weights are all equal, and so the variance of the weights, $\text{var}(w(y))$, is zero. Hence, we would like to choose the value of β which minimises the variance of the weights,

$$\beta_{\text{var}}^* = \arg \min_{\beta} \text{var}(w(y)).$$

In our experience, the mean of the estimator $\text{var}(w(y))$ is a smooth enough function of β that it can be used to tune the proposal variance during the burn-in phase of MCMC algorithms. However the variance of the estimator of the variance can be large, especially far away from the optimal value, so it can take a large number of iterations to calculate descent directions.

4.1.2. The effective sample size. The effective sample size, n_{eff} , can also 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 practise, we cannot achieve a perfect effective sample size of M .

The effective sample size can be 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 we see that when, if the variance of the weights is zero , then $n_{\text{eff}} = M$; this is our ideal scenario. Maximising the 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 varies between 1 and M as opposed to the variance which takes values on $[1, M]$ while the variance of the weights takes values on \mathbb{R}_+ . Moreover the variance of the weights can vary over many orders of magnitude . Therefore it is preferable as a means of tuning the scaling parameter. In all of the numerics which follow, we use causing numerical issues, 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 calculated at each iteration and averaged across the entire run to tune the scaling parameters. We do this because when we tune this parameter on the fly, we use the single-iteration average to estimate optimality. The optimal scaling parameter found using the global optimum of the variance of the weights is also included for comparison. Note that for the majority of iterations the value of 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} is an overestimate of this statistic over a larger number of samples. Rare events which result in a large importance weight bring this statistic down, and care must be made not to overfit the proposal variance to with scaling parameter δ_k . Calculating n_{eff} on an iteration by iteration basis. This can be accounted for by increasing the variance slightly once the adaptive algorithm has arrived on a value using single iteration values of 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 ~~imagine consider~~ the algorithm in the ~~burn-in phase~~ early stages, for example, we have ~~all M ensemble members in the tail of a Gaussian curve searching for the area of high density. If the ensemble members are evenly spaced, then the particles in the tails of the target searching for a mode. The~~ particle closest to the ~~mean mode~~ will have an exponentially higher weight assigned to it. ~~The and the~~ effective sample size ratio in this scenario will be close to 1. ~~As the algorithm burns in, the ensemble populates the regions where the majority of the probability density lies, and the proposal distributions better represent~~ In later stages the ensemble populates high density regions, and better represents the posterior distribution. This leads to smaller ~~variance variation~~ in the weights, and a ~~bigger higher~~ effective sample size. By ~~this argument we can see that rising n_{eff} signals the end of the burn-in period~~ looking for approximations of the effective sample size which look like a stationary distribution, we can tell when the algorithm is working efficiently.

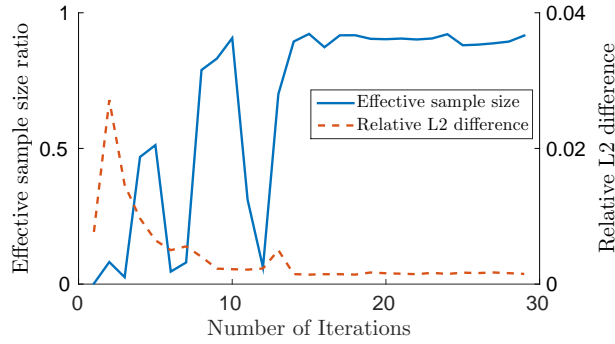
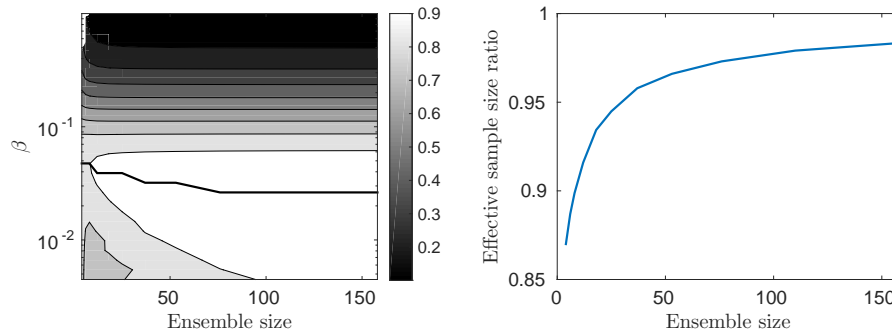


FIG. 4.1. The effective sample size ratio and relative L^2 difference ~~between E during the proposal and posterior distributions at each of the first 30 iterations. These numerics are taken from a simulation of the problem example in Section ??~~ 7.4.3 using the PAIS-Gamma-PAIS-RW algorithm.

Figure 4.1 ~~shows demonstrates~~ that the effective sample size flattens out ~~at the same time~~ as the relative L^2 difference between the posterior distribution and the proposal distribution stabilises close to its minimum. ~~Detecting this stationarity allows us to automatically determine the end of the burn-in phase of the algorithm.~~

~~If we are to use~~ Figure 4.2 shows how the effective sample size ratio as an indicator of how to tune the proposal variance, we need to look at how it behaves in different situations. Here we look at how the statistic behaves as we vary the ensemble size, M . Figure 4.2 shows results for the Gaussian posterior discussed in Section 7.2 when using the PAIS algorithm with RW proposals, each with variance $\beta^2 \in \mathbb{R}_{>0}$. Figure 4.2 ~~used~~ in PAIS is affected by the ensemble size. We see from subfigure (a) ~~shows that~~ as the ensemble size increases, the ~~scaling parameter which gives the optimal effective sample size ratio~~ optimal scaling parameter decreases. This is ~~to be expected since~~ if we are trying to approximate the posterior with a mixture distribution of a small number of Gaussians, the optimal variance will naturally be larger so that the whole of the significant regions are covered by the proposal distribution. As the number increases, the optimal variances decrease so that finer details in the posterior can be better represented in the proposal. Figure 4.2 ~~expected since the larger ensemble~~



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

FIG. 4.2. The behaviour of the effective sample size as the ensemble size increases—~~This analysis is from, considering the Gaussian posterior example in Section 7.2 using the PAIS-RW algorithm.~~

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 ~~efficiency of the sampler also increases~~ algorithm becomes more efficient.

4.2. Adaptive Adaptively Tuned PAIS. A popular approach for adaptive MCMC algorithms is to view ~~To adapt~~ the scaling parameter as a random variable which we can sample during the course of the MCMC iterations. However, it can be slow to converge to the optimal value, and we may need an uninformative prior for the scaling parameter. Alternatively, the parameter may be randomly sampled at various points during the evolution of the chain. This results in some iterations which make larger global moves in the state space between modes in the target distribution, and some which make local moves. Algorithms of this type do not converge to an optimal value of the scaling parameter. We choose to use a divide and conquer scheme which optimises the effective sample size (or any other diagnostic) β , we use a version of the gradient ascent method modified for a stochastic function. Some more sophisticated examples are described in [38] and [24]. [1, 24, 38]. From here on in, we will use the random walk proposal,

$$y = x_{i-1} + \beta \omega_{i-1}, \quad \omega_{i-1} \sim \mathcal{N}(0, \Sigma) \quad \text{i.i.d.}$$

where x_{i-1} is our current state, y is our proposed state and Σ is a covariance operator, which could come from the prior distribution. In the numerics section we will refer to β as the scaling parameter.

Using an adaptive strategy, we calculate a sequence $\{\beta^{(k)}\}_{k=1}$ which converges roughly to the optimal scaling parameter β^* , resulting in the optimal transition density χ for our MCMC algorithm. This optimal value will differ depending on the criterion we are optimising. We must choose some sequence of iterations, $\{n_k\}_{k=1}$, at which to update β , and due to the constraints on adaptive MCMC algorithms [37, 38], these. 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 5, the gradient ascent parameter γ may

Algorithm 2: Adaptively tuned PAIS algorithm.

```

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


```

~~decrease over time, e.g. as a function of n_k must grow exponentially further apart. This same adaptive approach can also be applied to the trivially parallelised MCMC algorithms to adaptively calculate their optimal scaling parameter β^* .~~

5. Approximate Multinomial Resampling. Although the ETPF is optimal in terms of preserving statistics of the sample, 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 (AMR).

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) algorithm.

```

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

```

Table ?? 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.

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. 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. 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.

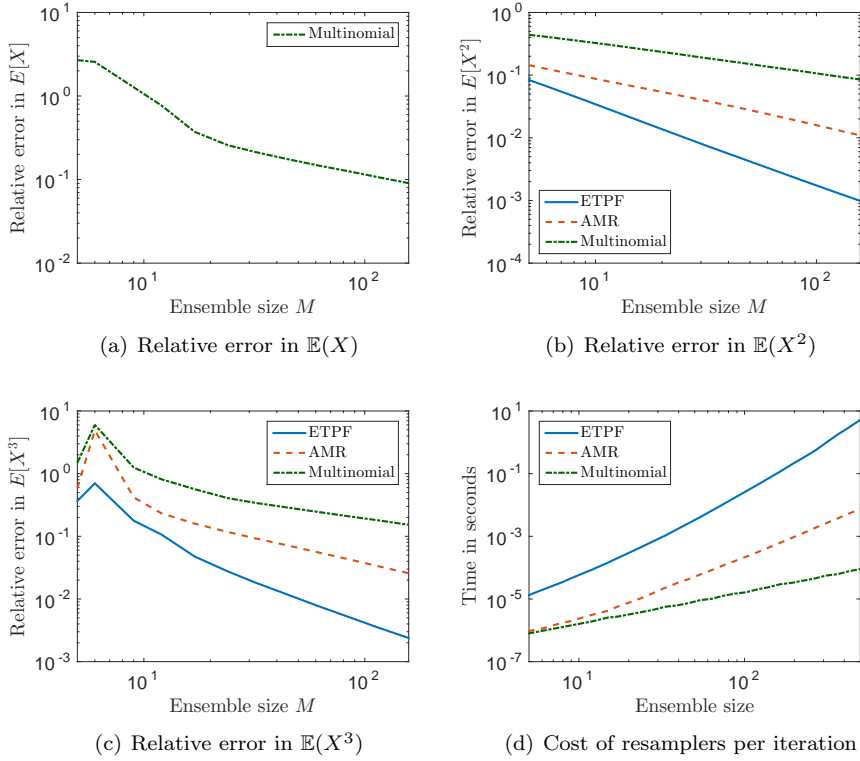


FIG. 5.1. ~~Finding optimal values~~ Comparison of β for the ~~problem in Section 7.2~~ performance between different resampling schemes. The ~~setup is as example in Section 7.1~~ β is implemented for this demonstration.

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 of the PAIS algorithm.

6. Numerical Examples.

5.1. Sampling from a one dimensional Gaussian distribution.

6. Consistency of PAIS. In this example we compare the naively parallelised RWMH algorithm with its PAIS variant, the PAIS-RW algorithm. The PAIS algorithm is implemented using the ETPF to perform the resampling step. We assess the performance of the PAIS algorithm using the relative L^2 error defined in β , as well as the relative error in the first moment. As outlined in [30] 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 $N \rightarrow \infty$.

Since we are comparing against naively parallelised MH algorithms, In case one, from standard importance sampling results we know that for an iteration n , as $M \rightarrow \infty$, we also need to decide which statistics $T(\beta)$ provide the best criterions for obtaining the optimal scaling parameters. In the examples which follow, we have optimised 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, but this higher acceptance rate is commonly used for one dimensional Gaussian posteriors [39]. To find the optimal scaling parameter we minimise the statistic obtain a consistent estimator for any statistic of interest.

$$\underline{T_{\text{MH}} \hat{\phi}(\beta \Theta) = \frac{N_{\text{acc}}(\beta)}{N_{\text{total}}} \approx \frac{1}{\hat{Z}} \sum_{i=1}^M \frac{1}{M} w_i \phi(\theta_i) \rightarrow \phi(\Theta),}$$

where $N_{\text{acc}}(\beta)$ is the number of accepted moves and N_{total} is the total number of samples produced. For the PAIS algorithm, we maximise the effective sample size as discussed in Section 4.1.2 the normalisation constant, $\hat{Z} = \frac{1}{M} \sum_{i=1}^M w_i$, also converges to the true normalisation constant Z [35].

6.0.1. Target distribution. Consider the simple case of a linear observation operator $\mathcal{G}(x) = x$, where the prior on x and the observational noise follow Gaussian distributions. Then, following , the Gaussian posterior has the form

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

where σ^2 and τ^2 are the variances of the observational noise and prior distributions respectively. In the numerics which follow, we choose $\tau^2 = 0.01$ Case two is slightly more involved. Estimation of the normalisation constant Z is biased, and $\sigma^2 = 0.01$, and we observe $x_{\text{ref}} = 4$ noisily such that

$$\underline{D = \mathcal{G}(x) + \eta \sim \mathcal{N}(\mathcal{G}(x_{\text{ref}}), \sigma^2) = \mathcal{N}(4, 0.01).}$$

These values result in a posterior density in which the vast majority of the density is out in so estimates of statistics are sums of independent but biased estimators. Since the estimators are independent, 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. 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 $\hat{Z} \rightarrow Z$ [35]. Since the normalisation constant is consistent, sums of the independent estimators are also consistent.

6.0.1. Numerical implementation. In each of the following simulations, we perform three tasks. First we calculate the optimal value of β by optimising the statistics described in Section 4.1. We then run the algorithms with optimal parameters to calculate and compare the convergence rates. Finally, we implement the adaptive algorithms described in Section 4.2 and compare the convergence rates of these algorithms with the nonadaptive algorithms

7. Numerical Examples. 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. ~~(1) Finding the optimal scaling parameters:~~ ~~To find the optimal parameters we choose~~ We chose 32 values of β evenly spaced on a log scale in the interval $[10^{-5}, 2]$. We ~~run the PAIS-RW and RWMH~~ ran the PAIS and MH algorithms for one million iterations, each with an ensemble size of $M = 50$. We took 32 repeats ~~of both algorithms and then and~~ used the geometric means of the sample statistics in Section 4.1 to find the optimal ~~parameters.~~ scaling parameter.
2. ~~(2) Measuring convergence of nonadaptive algorithms:~~ We ~~run the algorithms in Section 7.2~~ 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 ~~using the optimal parameters found in (1).~~ The performance ~~of the algorithms is judged by the convergence of the~~ is evaluated using the relative L^2 error ~~statistic in defined in Equation (4.2).~~
3. ~~(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 ~~$\beta^{(1)}$~~ $\beta_1 = 1$.

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 [39]. 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.

7.2.2. ~~Optimal values of β~~ 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 curve we can see that the minimum is very wide and despite the

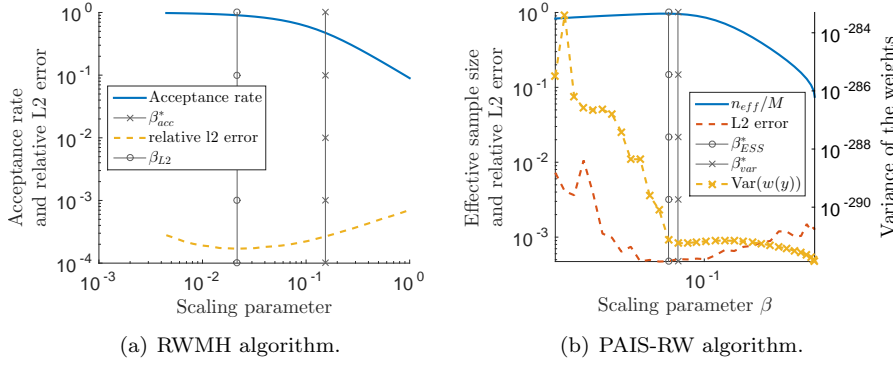


FIG. 7.1. Finding optimal values of β for the ~~problem example~~ in Section 7.2. The setup is as in Section 7.1. Resampling is performed using the ETPF.

Statistic	RWMH
β_{L2}^*	2.1e-2
$\beta_{\%}^*$	1.5e-1
Acceptance Rate (β_{L2}^*)	9.0e-1
Acceptance Rate ($\beta_{\%}^*$)	5.0e-1

Statistic	PAIS-RW
β_{eff}^*	4.7e-2
$\beta_{var(w(y))}^*$	5.8e-2
β_{L2}^*	3.9e-2

TABLE 7.1

Optimal values of β summarised from Figure 7.1. Statistics calculated as described in Section 4.1. The values β_{L2}^* and $\beta_{\%}^*$ are the optimal scaling parameters found by optimising the relative L^2 errors and acceptance rate respectively. Similarly β_{eff}^* and $\beta_{var(w(y))}^*$ optimise the effective sample size and variance of the weights statistics.

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. ~~In general this indicator overestimates the value of n_{eff} found by using the entire sample.~~

7.2.3. Convergence of RWMH vs PAIS-RW. Figure 7.2 shows that the PAIS-RW algorithm converges ~~to the 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.

~~Both The~~ adaptive algorithms are ~~run with initial values of $\beta = 1$. also shown in~~ Figure 7.2 ~~shows that after an initial burn-in period the APAIS-RW algorithm catches up to the PAIS-RW algorithm, and by the end of the simulation window is matching its performance. The ARWMH algorithm does not perform quite as well, this is possibly due to the fact that the acceptance rate cost function is not particularly smooth at the optimal value making it difficult to minimise. 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~~

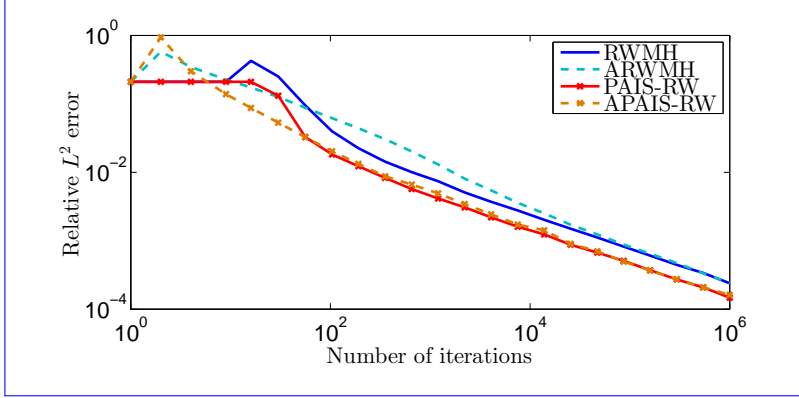


FIG. 7.2. ~~Error analysis for~~ Convergence of the (A)RWMH and (A)PAIS-RW algorithms. The setup is as in Section 7.1 (2, 3). Resampling is performed using the ETPF.

fly.

7.2.4. Scaling of the PAIS algorithm with ensemble size. Throughout ~~this example,~~ In this and the next example we use an ensemble size $M = 50$, ~~but~~ however it is interesting to see how the PAIS algorithm scales when we increase the ensemble size, and if there is ~~some limit below which the algorithm fails~~ 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 $M \in [1, 160]$.

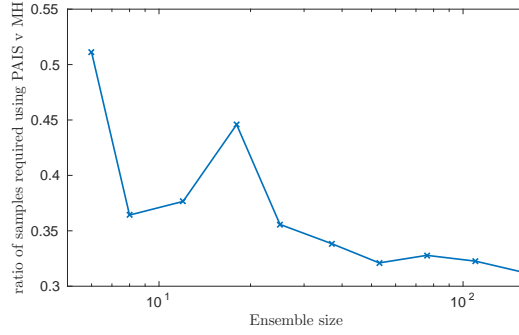


FIG. 7.3. *Ratio of PAIS-RW samples required to reach the same tolerance as the RWMH algorithm.*

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 ~~indicates~~ 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 ~~linked~~ due

to the increasing effective sample size shown in Figure 4.2 (b).

7.3. Sampling from Bimodal Distributions. In this ~~section~~ ~~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 ~~chains~~ ~~particles~~ to new modes as they are found. This means that we expect the number of ~~chains~~ ~~particles~~ in a mode to be approximately proportional to the probability density in that mode. ~~As a result, reconstructed posteriors with disproportional modes, as is familiar with the MH algorithms, are not produced, resulting in faster global convergence.~~

7.3.1. Target Distribution We look at an ‘easy’ problem, B_1 , which has a KL-divergence of 0.880, and a ‘harder’ problem, B_2 ,

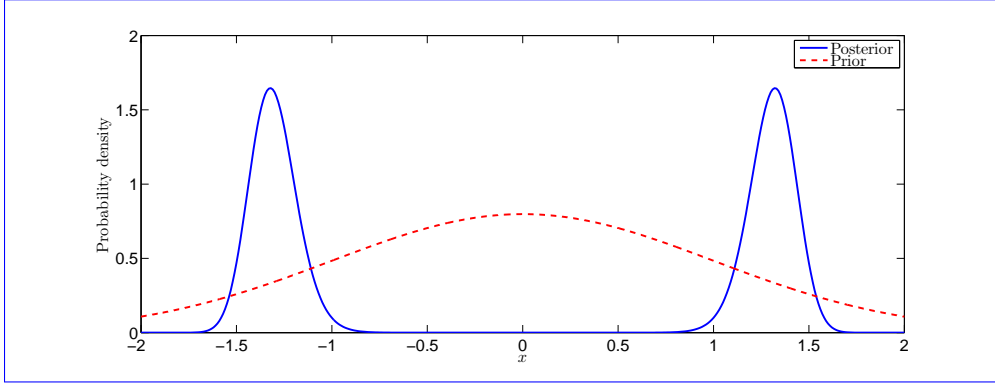


FIG. 7.4. *Posterior distribution for the bimodal example in Section 7.3.1.*

In this example we look at a bimodal posterior distribution which has a KL-divergence of 3.647. Problem B_1 has two modes which are not too far apart. In B_2 we increase the distance between the two modes which has the effect of increasing the high energy barrier between the two modes. In the MH algorithm, the expected number of iterations that it takes for a MCMC chain to jump between modes. These posteriors are shown for a switch between modes is a large proportion of the total number of iterations. This posterior is presented in Figure ??7.4.

The following setup is the same for both problems. We consider a non-linear observation operator $\mathcal{G}(x) = x^2$, and assign the prior $x \sim \mu_0 = \mathcal{N}(0, \tau^2 = 0.25)$. We assume that a noisy reading, D , is taken according to $D = \mathcal{G}(x_{\text{ref}}) + \epsilon$, where $\epsilon \sim \mu_\epsilon = \mathcal{N}(0, \sigma^2 = 0.1)$. This results in the non-Gaussian posterior $\mu_0 = \mathcal{N}(0, \tau^2 = 0.25)$. 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

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

To create the ‘easy’ problem we say that the true value of $\mathcal{G}(x_{\text{ref}}) = 0.75$, and the ‘hard’ problem is generated using $\mathcal{G}(x_{\text{ref}}) = 2$. In the numerics which follow we draw noise from μ_ϵ to generate our data point.

Posterior distributions for problems B_1 and B_2 as described in Section 7.3.1.

7.3.2. Calculating values of ~~Optimal β^*~~ the optimal scaling parameter.

Calculating the optimal values of the scaling parameters for this ~~problem example~~ is similar to the previous example; ~~we check only the acceptance rate to find the optimal values for RWMH and we use the effective sample size to find the optimal values for PAIS-RW.~~ Table ?? gives the optimal values of β for both problems. 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}^* RWMH 4.8e-1 PAIS-RW -~~

	Algorithm	β_{acc}^*	β_{eff}^*	β_{L2}^*
1.0e-1	RWMH	2.3e-1	-	9.3e-1
	PAIS-RW	-	5.1e-2	1.3e-1

TABLE 7.2

~~Optimal values of β scaling parameters for B_1 (left) and B_2 (right)~~ the example in Section 7.3.

It is relatively simple to find optimal scaling parameters for problem B_1 . These values are given in Table ?? (left). However problem B_2 is much harder as Since transitions between the modes are extremely unlikely ~~for the standard RWMH algorithm.~~ This means that we need to, we consider the convergence on two levels; ~~we should consider the algorithm's ability to find both the~~ convergence to equally proportioned modes, and ~~also whether it can sample them in the correct proportions.~~

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 ~~means that β must be large~~ 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 ~~by 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.

~~From Table ?? (right) we see that there is a large difference between the optimal value of β for each regime in RWMH, and so will result in inefficient sampling. The PAIS-RW algorithm manages to sample the local detail and the large scale behaviour with similar values of β^* ; a clear advantage to using this algorithm for this problem.~~

7.3.3. Convergence of RWMH vs PAIS-RW. For B_2 ~~the~~ The MH and PAIS algorithms are run with ~~the global optimal value of β^* , and with the local optimal value of β^*~~ both the globally and locally optimal scaling parameters. Figure 7.5 (a) shows that the ~~algorithms using the globally optimal β^* converge at the desired rate, whereas the algorithms optimised using the acceptance rate and effective sample size leads to normal convergence rates, $\mathcal{O}(n^{-1/2})$, whereas algorithms using locally optimal scaling parameters initially converge faster but at some point forget the location of one of the modes, causing the convergence to flatten out~~ do not sample proportionally for the whole simulation.

Using the PAIS algorithm we can ~~minimise the impact of forgotten modes by constantly allowing the algorithm to search them out~~ ensure that we do not lose modes as simulation progresses. Since we have parallel chains an ensemble, we can ~~run PAIS-RW with the majority of chains using the locally optimal value of β , and sacrifice locally~~

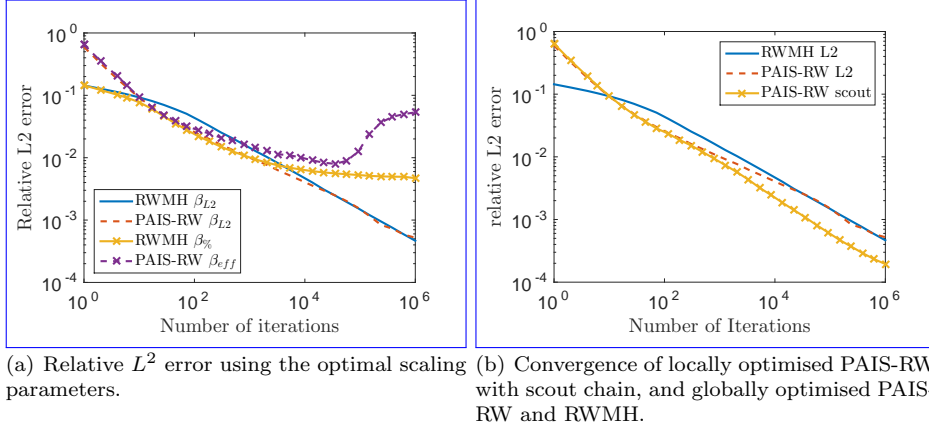


FIG. 7.5. Convergence of the PAIS-RW and RWMH algorithms for ~~Problem B_1~~ the bimodal example discussed in Section 7.3. Set up described in Section 7.1. Resampling is performed using the ETPF.

As in the Gaussian example we see a significant speed up with the PAIS-RW algorithm for problem B_1 . Figure ?? shows the adaptive and nonadaptive convergence rates. We can see that the adaptive algorithms compare closely with the respective nonadaptive algorithms and the improvement PAIS offers remains significant.—

Convergence of the PAIS-RW and RWMH algorithms for Problem B_2 . Set up described in Section 7.1. Resampling is performed using the ETPF.

optimal sampling for one or two chains with a larger scaling parameter 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 ~~other problems. Other methods of mode searches are described in [26].~~Gaussian example.

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 with the same proposal distributions, and each group will also have a ensemble members tuned for locally optimal convergence and one ‘scout’ chain with a scaling parameter ten times that of the rest of the group.

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.

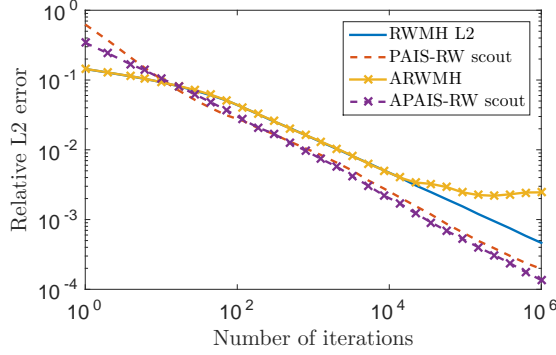


FIG. 7.6. Convergence of ~~the relative L^2 error for problem B2~~ modified PAIS and RWMH algorithms, comparing the globally optimised nonadaptive algorithms with the locally optimised adaptive algorithms. The setup is as described in Section 7.1 (3). Resampling is performed using the ETPF.

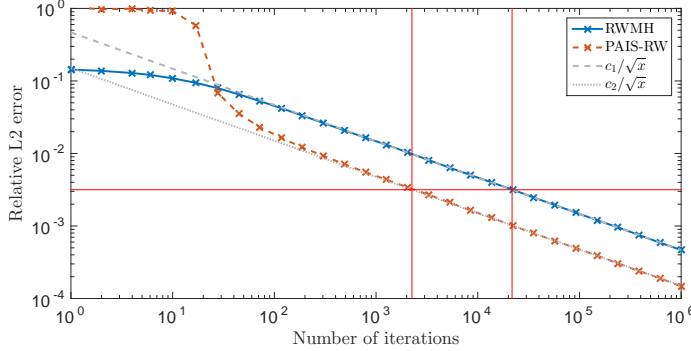


FIG. 7.7. Illustration of calculating the number of ~~PAIS-pCNL-PAIS-RW~~ iterations required to reach a tolerance of ~~10^-2.8~~ 10^-2.8 as a percentage of MH iterations. The relative L^2 error graphs are from the Gaussian problem in Section 7.2.

7.3.4. Calculating the ~~Speed Up~~ speed up in Convergence.

The graphs in the previous section clearly show that the PAIS-RW algorithm converges faster than the RWMH 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 ~~graph~~ algorithm has not finished burning in.

~~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. The blank entries correspond to occasions when either the MH algorithm or PAIS algorithm hasn't converged to the posterior distribution.~~

7.3.5. ~~A Useful Property of the PAIS Algorithm for Multimodal Distributions.~~

~~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.~~

	Gaussian	B_1 - B_2 -Bimodal
RWMH	10%	32% 12.6% (scout)
pCN 32%—MALA	42%	36% 40%
pCNL 66% 56%—		

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. ~~The pCN and pCNL proposal distributions are taken from [14].~~

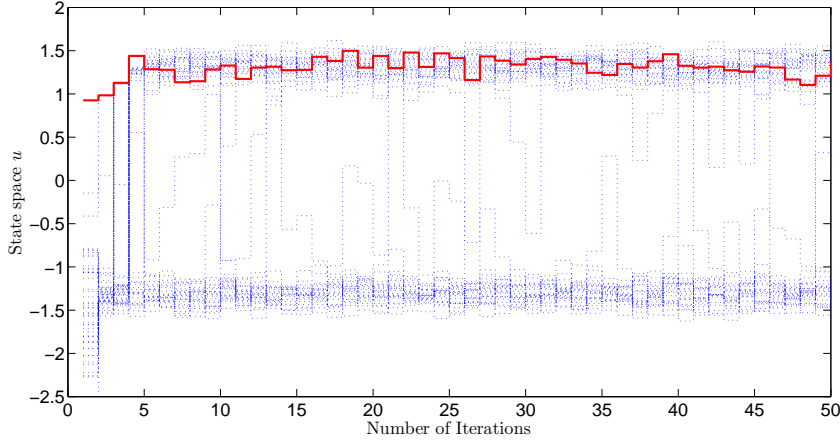


FIG. 7.8. This figure demonstrates the redistribution property of the PAIS algorithm. Initially there is one chain in the positive mode, and 49 chains in the negative mode.

7.3.5. A useful property of the PAIS algorithm for multimodal distributions.

The biggest issue for the Metropolis-Hastings algorithms when sampling from ~~a posterior~~ such as the one in B_2 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 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 ~~Problem B_2~~ 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 middle.

7.4. ~~Sampling from non-Gaussian bivariate distributions~~ A Mixture Model.

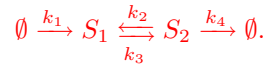
~~In this section we apply the PAIS algorithm to a more complicated posterior distribution. The field of biochemical kinetics gives rise to multiscale stochastic problems which remain a challenge both theoretically and computationally. Biochemical reactions~~

occur in single cells between a number of chemical populations and the rates of these reactions can vary on vastly different timescales. It is these reaction rates which we are interested in finding descriptions for.

It is often possible to isolate which reactions are occurring more frequently (the fast reactions) and which are occurring less frequently (the slow reactions). The quasi-steady-state assumption (QSSA) is the assumption that the fast reactions converge in distribution on a timescale which is negligible with respect to

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 rate of occurrence of the slow reactions. This assumption allows us to approximate the dynamics of the slowly changing quantities in the system by assuming that the fast quantities are in equilibrium with respect to the fast reactions in isolation. This kind of model reduction can be used to approximate the likelihood in an inverse problem where we wish to recover the reaction parameters in the system. 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.

Let us consider a simple example by introducing the following simple chemical system involving two chemical species S_1 and S_2 :



Each arrow represents a reaction from a reactant to a product, with some rate constant k_i , and where the rates of the reactions are assumed to follow mass action kinetics. We denote the concentration of species S_i by X_i . We assume that we are in a parameter regime such that the reactions $S_1 \rightarrow S_2$ and $S_2 \rightarrow S_1$ occur much much more frequently than the other reactions. Notice that both chemical species are involved in fast reactions. However, the quantity $\mathcal{S} = X_1 + X_2$ is conserved by both of the fast reactions, and as such, this is the slowly changing quantity in this system. The effective dynamics of \mathcal{S} can be represented as follows.



Here, the new reaction rate \hat{k}_4 is approximated through application of the QSSA to be

$$\hat{k}_4(s) = \mathbb{E}[k_4 X_2 | \mathcal{S} = s] = \frac{k_2 k_4 \mathcal{S}}{k_2 + k_3}.$$

The value of $\mathbb{E}[k_4 X_2 | \mathcal{S} = s]$ is approximated by finding the steady state of the ODE representing the fast subsystem of reactions:



If we assume that we know the rate constants k_1 and k_4 , and observe the system in \mathcal{S} , we indirectly observe the rates k_2 and k_3 through the effective rate \hat{k}_4 of the degradation of \mathcal{S} . Our observations are uninformative about these reaction rates, as there are surfaces in parameter space along which the effective rate \hat{k}_4 is invariant, leading to a highly ill-posed inverse problem. Making the assumption that the errors in our observations of the value of \mathcal{S} are Gamma distributed with some variance σ^2 , this results in a long and thin posterior distribution on k_2 and k_3 . This type of problem is notoriously difficult to sample from using 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, as the algorithms quickly find a point on this manifold on which \hat{k}_4 is invariant, but exploration along its length is slow 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. We now formalise the posterior of interest. We look for a distribution over the parameter $\mathbf{k} = (k_2, k_3)^T$, given that we know $k_1 = 100$. 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 $k_4 = 1$. We generate our data by making ten observations of the system at $t_i = 2, 4, \dots, 20$, here simulated by solving the full system, governed by the differential equations—

$$\begin{aligned}\frac{dS_1}{dt} &= k_1 - k_2 X_1(t) + k_3 X_2(t), \\ \frac{dS_2}{dt} &= k_2 X_1(t) - (k_3 + k_4) X_2(t),\end{aligned}$$

with initial conditions $X_1(0) = X_2(0) = 0$, and parameter values $\mathbf{k} = (50, 100)^T$. We then add noise taken from a Gamma distribution with variance $\sigma^2 = 225$ and centred at each $\mathcal{S}(t_i) = X_1(t_i) + X_2(t_i)$.

In our modelling we use the QSSA to simplify this system of differential equations into the one dimensional system based on—

$$\frac{d\mathcal{S}}{dt} = k_1 - \frac{k_2 k_4}{k_2 + k_3} \mathcal{S}(t).$$

The observation operator, $\mathcal{G} : (\mathbf{k}, t) \mapsto \mathcal{S}(t)$, maps from parameter space onto population space at a time t . This means that for the i th observation we assume 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 $p \in [0, 1]$,

$$D_i \sim \text{Gamma}(\alpha_i, \beta_i),$$

$\mu_1, \mu_2 \in \mathbb{R}$ 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

$$\alpha_i = \frac{\mathcal{G}(\mathbf{k}, t_i)^2}{\sigma^2}, \quad \beta_i = \frac{\mathcal{G}(\mathbf{k}, t_i)}{\sigma^2}, \quad i = 1, \dots, 10.$$

These are parameters required to give us a distribution with mean $\mathcal{G}(\mathbf{k})$ and variance σ^2 . We assign Gamma priors to \mathbf{k} with mean $\alpha_0/\beta_0 = 75$ and variance $\alpha_0/\beta_0^2 = 100$ in both coordinates, resulting in the posterior

$$\pi(\mathbf{k}|\mathbf{D}) \propto \left[\prod_{i=1}^{10} \text{Gamma}(D_i; \alpha_i, \beta_i) \right] \text{Gamma}(k_2; \alpha_0, \beta_0) \text{Gamma}(k_3; \alpha_0, \beta_0).$$

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

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

7.4.2. Implementation. For this problem there is

We have no analytic form for the normalisation constant, and numerical methods implemented in MATLAB and Mathematica have proven to be unreliable. To demonstrate the convergence of the algorithm, we first run a MH algorithm for much longer than we normally would (8×10^{10} samples), and consider the histogram produced to be a well converged approximation of the true posterior distribution. We then perform our usual simulations with the PAIS and MH algorithms to compare the rate at which the histograms converge to the posterior over the same mesh.

In this problem we modify our PAIS algorithm to use a mixture of Gamma distributions in the proposal distribution instead of a Gaussian mixture. The PAIS-Gamma 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 MH-Gamma algorithms use a Gamma proposal distribution with mean centred at the previous state,

$$y \sim \text{Gamma}(\cdot; \alpha^*, \beta^*) \quad \text{where} \quad \frac{\alpha^*}{\beta^*} = x \quad \text{and} \quad \frac{\alpha^*}{(\beta^*)^2} = \beta^2.$$

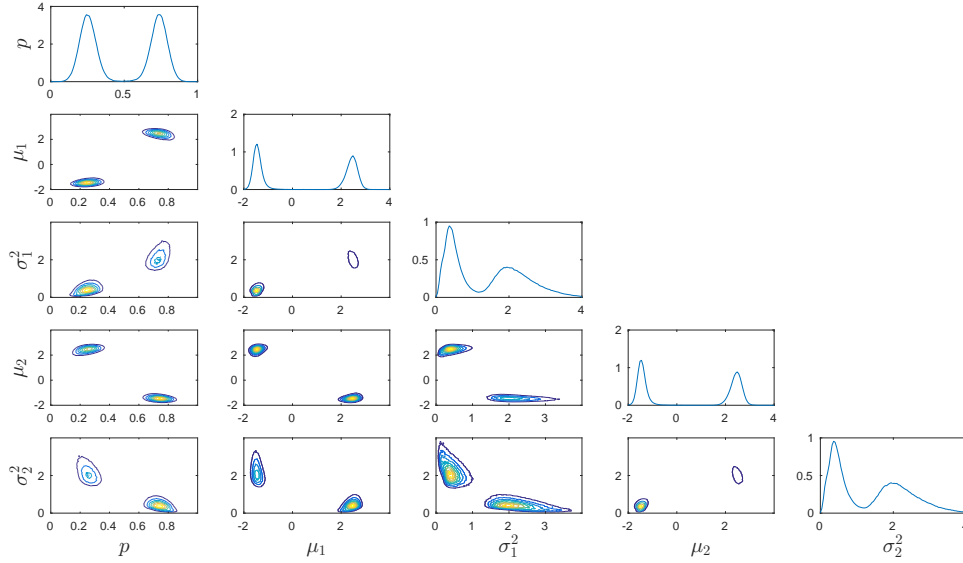


FIG. 7.9. The posterior distribution on the parameters k_2 and k_3 of θ as given in Equation (7.2) as found from 10 million samples from the data-PAIS algorithm. The main diagonal contains the marginal distributions of each θ_i and priors described in Section ?? the lower triangular contours represents the correlation between pairs of parameters.

Similarly the PAIS-LA algorithm uses the Langevin proposal distribution from the standard MALA algorithm with a perturbation taken from a Gamma distribution,

$$y \sim \text{Gamma}(\cdot; \alpha^*, \beta^*), \quad \text{where} \quad \frac{\alpha^*}{\beta^*} = x + \frac{1}{2} \beta^2 \nabla \log \pi(X) \quad \text{and} \quad \frac{\alpha^*}{(\beta^*)^2} = \beta^2.$$

This ensures that we only propose positive parameters. Having a heavier right tail in 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 proposal distribution also means that inverse problem i.e. the value of θ used to generate the data, and also the θ obtained by a relabelling of the posterior is absolutely continuous with respect to the proposal. When this is not the case, the weight function can tend to zero or infinity when $k_i \rightarrow \infty$, which results in poor, very spiky approximations of the tails of the posterior distribution leading to slow convergence. 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

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

We now significantly increase the ensemble size we are using from $M = 50$ to $M = 2500$. This allows us to build a better approximation of the posterior distribution for our proposals. Following the discussion in Section 5, it is clear that to keep the runtime of the resampler negligible compared with the calculation of the posterior we must switch from the ETPF resampler to the AMR algorithm. 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.

7.4.3. Convergence of PAIS-Gamma and PAIS-LA vs MH with Gamma proposals. In the numerics which follow we have used the AMR algorithm to resample with an ensemble size of $M = 2500$ 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 sections examples. We perform test runs to find the optimal scaling parameters for both Gamma proposals and MALA-type proposals considering convergence to modes with equal density. We then calculate the convergence rates of the algorithms by producing 50-10 million samples from the posterior with each algorithm, and repeat the simulation 32 times. Convergence of the PAIS-Gamma and PAIS-LA algorithms for the chemical system problem described in Section ?? . Implementation described in Sections ?? and ??. Resampling is performed using the AMR scheme.

Algorithm	MH-Gamma-MH	PAIS-Gamma-PAIS-LA-PAIS
$\beta^* \delta^*$	2.7e-0 9.41.3e-1	1.0e-0 2.3e-1

TABLE 7.4

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

7.4.3. Convergence of MH vs PAIS. In Figure ?? (a) we see that the two PAIS algorithms have similar optimal values of the scaling parameter, also shown in Table ??. However, the MALA-type proposal achieves a higher effective sample size with the same ensemble size. Convergence after the first 50 million samples is shown. 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.

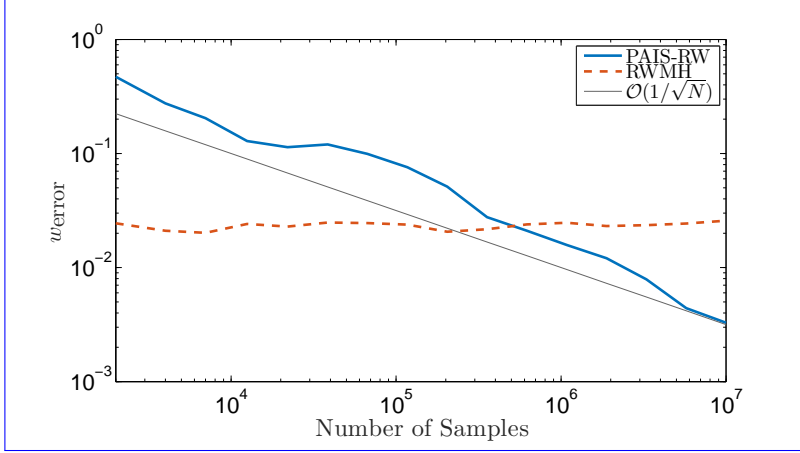


FIG. 7.10. *Convergence of the PAIS algorithm for the mixture model described in Section 7.4, convergence judged using the criterion in Equation (7.3). Implementation described in Sections 7.4.2 and 7.4.3. Resampling is performed using the AMR scheme.*

Convergence of the relative error for the two algorithms is displayed in Figure ?? (b) and demonstrates that the PAIS algorithm converges faster than the MH algorithm 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. The slightly unstable convergence in the PAIS Gamma algorithm is the effect of spikes appearing in the tails of the posterior distribution. This is also what has caused the slightly lower effective sample size.

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 [40].

Importantly, we have compared the efficiency of our parallel scheme with a naive 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.

Another strength of the PAIS is that it can also be used with any MCMC proposal. There are a growing number of increasing sophisticated MCMC algorithms (HMC 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 [33]) 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.

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] A. BESKOS, F. PINSKI, J. SANZ-SERNA, AND A. STUART, *Hybrid Monte Carlo on Hilbert spaces*, Stochastic Processes and their Applications, 121 (2011), pp. 2201–2230.
- [3] JORIS 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 pedigree 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] O. CAPPÉ, A. GUILLIN, J. MARIN, AND C. ROBERT, *Population Monte Carlo for ion channel restoration*.
- [9] O. CAPPÉ, A. GUILLIN, J. MARIN, AND C. ROBERT, *Population monte carlo*, Journal of Computational and Graphical Statistics, (2012).
- [10] G. CELEUX, J. MARIN, AND C. ROBERT, *Iterated importance sampling in missing data problems*, Computational statistics & data analysis, 50 (2006), pp. 3386–3404.
- [11] J. CORNUET, J. MARIN, A. MIRA, AND C. ROBERT, *Adaptive multiple importance sampling*, Scandinavian Journal of Statistics, 39 (2012), pp. 798–812.
- [12] 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.
- [13] 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.
- [14] 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.
- [15] T. EL MOSELHY AND Y. MARZOUK, *Bayesian inference with optimal maps*, Journal of Computational Physics, 231 (2012), pp. 7815–7850.
- [16] 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.
- [17] A. GELMAN AND D. B. RUBIN, *Inference from Iterative Simulation Using Multiple Sequences*, Statistical Science, 7 (1992), pp. 457–472.
- [18] 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.
- [19] N. GORDON, D. SALMOND, AND A. SMITH, *Novel approach to nonlinear/non-Gaussian Bayesian state estimation*, in IEE Proceedings F (Radar and Signal Processing), vol. 140, IET, 1993, pp. 107–113.
- [20] HEIKKI HAARIO, EERO SAKSMAN, AND JOHANNA TAMMINEN, *Componentwise adaptation for high dimensional mcmc*, Computational Statistics, 20 (2005), pp. 265–273.
- [21] W. HASTINGS, *Monte Carlo sampling methods using Markov chains and their applications*, Biometrika, 57 (1970), pp. 97–109.
- [22] THOMAS HOUSE, ASHLEY FORD, SHIWEI LAN, SAMUEL BILSON, ELIZABETH BUCKINGHAM-

- JEFFERY, AND MARK GIROLAMI, *Bayesian uncertainty quantification for transmissibility of influenza, norovirus and ebola using information geometry*, Journal of The Royal Society Interface, 13 (2016), p. 20160279.
- [23] M. ISARD AND A. BLAKE, *Condensation conditional density propagation for visual tracking*, International journal of computer vision, 29 (1998), pp. 5–28.
 - [24] 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.
 - [25] R. KALMAN, *A new approach to linear filtering and prediction problems*, Journal of Fluids Engineering, 82 (1960), pp. 35–45.
 - [26] S. LAN, J. STREETS, AND B. SHAHBABA, *Wormhole Hamiltonian Monte Carlo*, ArXiv e-prints, (2013).
 - [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] JEAN-MICHEL MARIN, KERRIE Mengersen, AND CHRISTIAN P ROBERT, *Bayesian modelling and inference on mixtures of distributions*, Handbook of statistics, 25 (2005), pp. 459–507.
 - [30] 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.
 - [31] G. MOORE ET AL., *Cramming more components onto integrated circuits*, Proceedings of the IEEE, 86 (1998), pp. 82–85.
 - [32] 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).
 - [33] 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.
 - [34] S. REICH, *A nonparametric ensemble transform method for Bayesian inference*, SIAM Journal on Scientific Computing, 35 (2013), pp. A2013–A2024.
 - [35] C. ROBERT AND G. CASELLA, *Monte Carlo statistical methods*, Springer Science & Business Media, 2013.
 - [36] G. ROBERTS AND J. ROSENTHAL, *Optimal scaling for various Metropolis-Hastings algorithms*, Statistical science, 16 (2001), pp. 351–367.
 - [37] ———, *Coupling and ergodicity of adaptive Markov chain Monte Carlo algorithms*, Journal of Applied Probability, 44 (2007), pp. 458–475.
 - [38] ———, *Examples of adaptive MCMC*, Journal of Computational and Graphical Statistics, 18 (2009), pp. 349–367.
 - [39] J. ROSENTHAL, *Optimal proposal distributions and adaptive MCMC*, Handbook of Markov Chain Monte Carlo, (2011), pp. 93–112.
 - [40] P. RUSSELL, *PhD Thesis*, PhD thesis, School of Mathematics, University of Manchester, 2017.
 - [41] J. SEXTON AND D. WEINGARTEN, *Hamiltonian evolution for the hybrid Monte Carlo algorithm*, Nuclear Physics B, 380 (1992), pp. 665–677.
 - [42] J. SIRÉN, P. MARTTINEN, AND J. CORANDER, *Reconstructing population histories from single nucleotide polymorphism data*, Molecular biology and evolution, 28 (2011), pp. 673–683.
 - [43] ANTTI SOLONEN, PIIRKKA OLLINAHO, MARKO LAINE, HEIKKI HAARIO, JOHANNA TAMMINEN, HEIKKI JÄRVINEN, ET AL., *Efficient mcmc for climate model parameter estimation: parallel adaptive chains and early rejection*, Bayesian Analysis, 7 (2012), pp. 715–736.
 - [44] A. STUART, *Inverse problems: a Bayesian perspective*, Acta Numerica, 19 (2010), pp. 451–559.
 - [45] C. VILLANI, *Topics in Optimal Transportation*, Graduate studies in mathematics, American Mathematical Society, 2003.
 - [46] ———, *Optimal Transport: Old and New*, Grundlehren der mathematischen Wissenschaften, Springer Berlin Heidelberg, 2008.