

# Subsampling Sequential Monte Carlo for Static Bayesian Models

David Gunawan<sup>1,4</sup>, Khue-Dung Dang<sup>2,4</sup>, Matias Quiroz<sup>2,4,5</sup>,  
Robert Kohn<sup>3,4</sup> and Minh-Ngoc Tran<sup>4,6</sup>

## Abstract

We show how to speed up Sequential Monte Carlo (SMC) for Bayesian inference in large data problems by data subsampling. SMC sequentially updates a cloud of particles through a sequence of distributions, beginning with a distribution that is easy to sample from such as the prior and ending with the posterior distribution. Each update of the particle cloud consists of three steps: reweighting, resampling, and moving. In the move step, each particle is moved using a Markov kernel; this is typically the most computationally expensive part, particularly when the dataset is large. It is crucial to have an efficient move step to ensure particle diversity. Our article makes two important contributions. First, in order to speed up the SMC computation, we use an approximately unbiased and efficient annealed likelihood estimator based on data subsampling. The subsampling approach is more memory efficient than the corresponding full data SMC, which is an advantage for parallel computation. Second, we use a Metropolis within Gibbs kernel with two conditional updates. A Hamiltonian Monte Carlo update makes distant moves for the model parameters, and a block pseudo-marginal proposal is used for the particles corresponding to the auxiliary variables for the data subsampling. We demonstrate both the usefulness and limitations of the methodology for estimating four generalized linear models and a generalized additive model with large datasets.

**Keywords.** Hamiltonian Monte Carlo, Large datasets, Likelihood annealing

---

<sup>1</sup>:School of Mathematics and Applied Statistics, University of Wollongong. <sup>2</sup>:School of Mathematical and Physical Sciences, University of Technology Sydney <sup>3</sup>:School of Economics, UNSW Business School, University of New South Wales. <sup>4</sup>:ARC Centre of Excellence for Mathematical and Statistical Frontiers (ACEMS). <sup>5</sup>:Research Division, Sveriges Riksbank. <sup>6</sup>:Discipline of Business Analytics, University of Sydney.

# 1 Introduction

The aim of Bayesian inference is to obtain the posterior distribution of unknown parameters, and in particular the posterior expectations of functions of the parameters. This is usually done by obtaining a simulation approximation of the expectation using samples from the posterior distribution. Exact approaches such as Markov Chain Monte Carlo (MCMC) (Brooks et al., 2011) have been the main methods used for sampling from complex posterior distributions. Despite this, MCMC methods have some notable drawbacks and limitations. One drawback, often overlooked by practitioners when fitting complex models, is the failure to converge caused by poorly mixing chains. While Hamiltonian Monte Carlo (Neal, 2011, HMC) is a remedy in many cases, it can be notoriously difficult to tune. Limitations of MCMC methods include the difficulties of assessing convergence, parallelizing the computation, and estimating the marginal likelihood efficiently from MCMC output, the latter being useful for model selection (Kass and Raftery, 1995). Sequential Monte Carlo (see Doucet et al., 2001 for an introductory overview) methods provide an alternative exact simulation approach to MCMC methods and overcome some of their drawbacks. Moreover, in contrast to MCMC methods, SMC can provide online updates of the parameters as data is collected, which is particularly useful for dynamic (time-varying parameters) models. SMC is also useful for static (non time-varying parameters) models (Chopin, 2002; Del Moral et al., 2006), and can in such cases more easily explore multimodal posterior distributions than MCMC. Note that our definition of dynamic refers to the model parameters or any unobserved states being time-varying and not the data. For example, an autoregressive (AR) model is considered to be static as the parameters do not depend on time, whereas a state space model is considered to be dynamic since the states evolve through time.

Despite the advantages of SMC, it is remarkably less used than MCMC for static models. One possible explanation is that, while amenable to computer parallelization, it is still very computationally expensive and particularly so for large datasets. Another obstacle caused by large datasets is that they prevent efficient computer parallelization of SMC, as the full dataset needs to be available for each worker which is infeasible as it consumes too much Random-Access Memory (RAM). We propose an efficient data subsampling approach which significantly reduces both the computational cost of the algorithm and the memory requirements when parallelizing: see Section 3.6 for a detailed explanation of the latter. Our approach utilizes the methods previously developed for Subsampling MCMC (Quiroz et al., 2019; Dang et al., 2019) and places them within the SMC framework. See Quiroz et al. (2018b) for an introduction to Subsampling MCMC.

In the Bayesian context, SMC traverses a cloud of particles through a sequence of distributions, with the initial distribution both easy to sample from and to evaluate, while the final distribution is the posterior distribution. The cloud of particles at step  $p$  is an estimate of the  $p$ th distribution in the sequence. The particles consist of the unknown parameters and any additional latent variables that are part of the model. The evolution of the particle cloud from one step to another consists of three steps: reweighting, resampling and moving. Of these, the first two steps are common to all SMC schemes and are straightforward. The move step is the most expensive and is critical to ensure that the particle cloud is representative of the distribution it aims to estimate.

To the best of our knowledge, data subsampling has not been explored in SMC. While Wang et al. (2019) term their algorithm Subsampling SMC, their approach is distinct as they combine data annealing and likelihood annealing, whereas we use data subsampling to estimate the likelihood. In particular, data annealing requires handling all the data, whereas the data subsampling approach only deals with a small fraction of the data at each stage. Specifically, we consider a likelihood annealing approach in which we estimate the annealed likelihood efficiently using an approximately unbiased estimator. Likelihood estimates for SMC in a non-subsampling context have been used in Duan and Fulop (2015), who propose to estimate the likelihood unbiasedly using a particle filter in a time series state space model application. However, Duan and Fulop (2015) use a random walk MCMC kernel for the move step of the model parameters, which is inefficient in high dimensions and we now turn to this issue.

The literature has focused on accelerating SMC algorithms by designing efficient MCMC kernels for the move step to achieve efficient particle diversity. Efficiency here means the ability of the MCMC kernel to generate distant proposals which have a high probability of being accepted. The advantage of an efficient move step is that few iterations of the kernel are needed, which is computationally cheap. Various approaches exist to achieve this. For example, adaptive SMC adapts the tuning parameters of the kernel to improve its efficiency (Jasra et al., 2011; Fearnhead and Taylor, 2013; Buchholz et al., 2018). South et al. (2016) use SMC with a flexible copula based independent proposal, while Sim et al. (2012) and South et al. (2017) use derivatives to construct efficient proposals through the Metropolis Adjusted Langevin Algorithm (Roberts and Stramer, 2002, MALA). It is now well-known that the MALA proposal is a special case of the more general proposal utilizing Hamiltonian dynamics proposed in Duane et al. (1987) (see Neal (2011) and Betancourt (2017) for an introduction to HMC). Although South et al. (2017) mention HMC in their introduction, they only consider MALA in their paper and

show how neural networks can be applied to adaptively choose its tuning parameters. Daviet (2016) considers HMC proposals for particle diversity. However, HMC is very slow for very large datasets and therefore this approach does not scale well in the number of observations.

We propose data subsampling to achieve scalability in the number of observations and HMC Markov move steps to achieve particle diversity. Section 3.6 shows that data subsampling lowers the memory requirements of the algorithm, making it possible to parallelise the computing on very large datasets. Our framework combines that of Duan and Fulop (2015) for carrying out SMC with an estimated likelihood, Quiroz et al. (2019) for estimating the likelihood and controlling the error in the target density and Dang et al. (2019) for constructing efficient proposals for high-dimensional targets in a subsampling context.

The rest of the article is organized as follows. Section 2 reviews SMC for static models. Section 3 outlines the methodology. Section 4 applies the methodology in a variety of settings for simulated data. Section 5 presents an application of our method in model selection for a real dataset. Section 6 concludes.

## 2 Sequential Monte Carlo

### 2.1 SMC for static Bayesian models

Denote the observed data  $\mathbf{y} = (\mathbf{y}_1^\top, \dots, \mathbf{y}_n^\top)^\top$ , with  $\mathbf{y}_k \in \mathcal{Y} \subset \mathbb{R}^{d_y}$ , where  $\mathbb{R}^m$  is an  $m$  dimensional Euclidean space. Let  $\boldsymbol{\theta}$  be the vector of unknown parameters,  $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^{d_\theta}$ , with  $p(\boldsymbol{\theta})$  and  $p(\mathbf{y}|\boldsymbol{\theta})$  the prior and likelihood. In Bayesian inference, the uncertainty about  $\boldsymbol{\theta}$  is specified by the posterior density  $\pi(\boldsymbol{\theta})$ , which by Bayes' theorem is

$$\pi(\boldsymbol{\theta}) = \frac{p(\boldsymbol{\theta})p(\mathbf{y}|\boldsymbol{\theta})}{p(\mathbf{y})}, \quad (1)$$

where  $p(\mathbf{y}) = \int_{\Theta} p(\mathbf{y}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}$  is the marginal likelihood which is often used for Bayesian model selection.

An important problem in Bayesian inference is to estimate the posterior expectation of a function  $\varphi$  of  $\boldsymbol{\theta}$ ,

$$\mathbb{E}_\pi(\varphi(\boldsymbol{\theta})) = \int_{\Theta} \varphi(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}. \quad (2)$$

In simulation based inference, this is typically achieved by sampling from (1) and computing (2) by Monte Carlo integration. Another important problem is to compute the marginal likelihood in (1). However, it is well known that standard Monte Carlo

integration is very inefficient for this task.

SMC (Doucet et al., 2001; Del Moral et al., 2006) is a collection of methods that provide a convenient approach to computing the posterior distribution and in addition the marginal likelihood. Likelihood tempered SMC specifies a sequence of  $P$  densities, connecting the density of the prior  $p(\boldsymbol{\theta})$  to the density of the posterior  $\pi(\boldsymbol{\theta})$  in (1). The sequence is obtained through temperature annealing (Neal, 2001), in which the likelihood is tempered as  $p(\mathbf{y}|\boldsymbol{\theta})^{a_p}$  with  $a_0 = 0 < a_1 < \dots < a_P = 1$ . We note that frequently  $P$  as well as  $a_1, \dots, a_P$  are chosen adaptively as the SMC proceeds, and we do so in our article; see Section 2.2. Our article estimates the tempered likelihood  $p(\mathbf{y}|\boldsymbol{\theta})^{a_p}$  by data subsampling as in Section 3. The  $p$ th tempered posterior is

$$\pi_p(\boldsymbol{\theta}) = \frac{\eta_p(\boldsymbol{\theta})}{Z_p}, \text{ where } \eta_p(\boldsymbol{\theta}) = p(\mathbf{y}|\boldsymbol{\theta})^{a_p} p(\boldsymbol{\theta}) \quad \text{and} \quad Z_p = \int_{\boldsymbol{\Theta}} p(\mathbf{y}|\boldsymbol{\theta})^{a_p} p(\boldsymbol{\theta}) d\boldsymbol{\theta}. \quad (3)$$

SMC starts by sampling a set of  $M$  particles from the prior  $p(\boldsymbol{\theta})$  and traverses them through the sequence of densities  $\pi_p(\boldsymbol{\theta}), p = 1, \dots, P$  such that, for each  $p$ , the reweighting, resampling and move steps are performed on the particles. Here, we assume for simplicity that it is possible to sample from the prior; otherwise one can sample from some initial distribution  $\pi_0(\boldsymbol{\theta})$  whose support covers that of the prior  $p(\boldsymbol{\theta})$ . At the final  $p = P$ , the particles are a (weighted) sample from  $\pi(\boldsymbol{\theta})$ . We now discuss this in more detail.

The initial particle cloud and weights  $\{\boldsymbol{\theta}_{1:M}^{(0)}, W_{1:M}^{(0)}\}$  are obtained by generating the  $\{\boldsymbol{\theta}_{1:M}^{(0)}\}$  from  $p(\boldsymbol{\theta})$ , and giving them equal weight, i.e.,  $W_{1:M}^{(0)} = 1/M$ . The weighted particles  $\{\boldsymbol{\theta}_{1:M}^{(p-1)}, W_{1:M}^{(p-1)}\}$  at the  $(p-1)$ st stage,  $p = 1, \dots, P$ , are (weighted) samples from  $\pi_{p-1}(\boldsymbol{\theta})$ . At the  $p$ th stage, the transition from  $\pi_{p-1}(\boldsymbol{\theta})$  to  $\pi_p(\boldsymbol{\theta})$  is obtained by the *reweighting step*,

$$w_i^{(p)} = W_i^{(p-1)} \frac{\eta_p(\boldsymbol{\theta}_i^{(p-1)})}{\eta_{p-1}(\boldsymbol{\theta}_i^{(p-1)})} = W_i^{(p-1)} p(\mathbf{y}|\boldsymbol{\theta}_i^{(p-1)})^{a_p - a_{p-1}},$$

and then normalizing  $W_i^{(p)} = w_i^{(p)} / \sum_{i'=1}^M w_{i'}^{(p)}$ . The reweighting assigns vanishingly small weights to particles which are unlikely under the tempered likelihood. This might cause the so-called particle degeneracy problem, in which the weight mass is concentrated only on a small fraction of the particles, causing a small effective sample size (explained in Section 2.2). This is resolved by the *resampling step*, in which the particles  $\boldsymbol{\theta}_{1:M}^{(p)}$  are sampled with a probability equal to their normalized weights  $W_{1:M}^{(p)}$ , and then setting  $W_{1:M}^{(p)} = 1/M$ . We use multinomial resampling for all

the experiments and applications in the paper. While this ensures that the particles with small weights are eliminated, it causes the so-called particle depletion problem because resampling might lead to only a few distinct particles. This is resolved by the *move step*, in which a  $\pi_p$ -invariant Markov kernel  $K_p$  is applied to move each of the particles  $R$  steps. Since a particle after the resampling step at stage  $p$  is approximately a sample from  $\pi_p(\theta)$  and  $K_p$  is  $\pi_p$ -invariant, no burn-in period is required as in MCMC methods, where often a very large number of burn-in iterations are required. Finally, we note that the algorithm is easy to parallelize with respect to the  $M$  particles, because the computations required for each particle do not depend on those of the other particles. Thus, provided that  $p(\mathbf{y}|\theta)$  can be computed at each worker without storage issues, it is straightforward to implement the parallel version.

Del Moral et al. (2006) provide consistency results and central limit theorems for estimating (2) based on the SMC output.

## 2.2 Statistical efficiency of SMC

The statistical efficiency of the  $p$ th stage of the SMC reweighting part is measured through the Effective Sample Size (ESS) defined as (Liu, 2001)

$$\text{ESS}_p := \left( \sum_{i=1}^M \left( W_i^{(p)} \right)^2 \right)^{-1}.$$

The  $\text{ESS}_p$  varies between 1 and  $M$ , where a low value of  $\text{ESS}_p$  indicates that the weights are concentrated only on a few particles. It is necessary to choose the tempering sequence  $\{a_p, p = 1, \dots, P\}$  carefully because it has a substantial impact on the  $\text{ESS}_p$ . We follow Del Moral et al. (2012) and choose the tempering sequence adaptively to ensure a sufficient level of particle diversity by selecting the next value of  $a_p$  such that  $\text{ESS}_p$  stays close to some target value  $\text{ESS}_{\text{target}}$ ; this is done by evaluating the  $\text{ESS}_p$  over a grid points  $a_{1:S,p}$  of potential values of  $a_p$  for a given  $p$  and selecting  $a_p$  as that value of  $a_{s,p}$ ,  $s = 1, \dots, S$ , whose  $\text{ESS}_p$  is closest to  $\text{ESS}_{\text{target}}$ . Throughout our article  $\text{ESS}_{\text{target}} = 0.8M$ .

For this adaptive choice of tempering sequence, Beskos et al. (2016) establish consistency results and central limit theorems for estimating (2) based on the SMC output. Other adaptive methods to choose the tempering sequence such as the approach by Del Moral et al. (2012) may also be used instead.

## 2.3 SMC estimation of the marginal likelihood

The marginal likelihood  $p(\mathbf{y})$  is often used in the Bayesian literature to compare models by their posterior model probabilities (Kass and Raftery, 1995). An advantage of SMC is that it automatically produces an estimate of  $p(\mathbf{y})$ .

Using the notation of Section 2.1,  $Z_P = p(\mathbf{y})$ ,  $Z_0 = 1$ , and

$$p(\mathbf{y}) = \prod_{p=1}^P \frac{Z_p}{Z_{p-1}} \quad \text{with} \quad \frac{Z_p}{Z_{p-1}} = \int \left( \frac{\eta_p(\boldsymbol{\theta})}{\eta_{p-1}(\boldsymbol{\theta})} \right) \pi_{p-1}(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$

Because the particle cloud  $\{\boldsymbol{\theta}_{1:M}^{(p-1)}, W_{1:M}^{(p-1)}\}$  at the  $(p-1)$ st stage is an approximate sample from  $\pi_{p-1}(\boldsymbol{\theta})$ , the ratios above are estimated by

$$\frac{\widehat{Z}_p}{Z_{p-1}} = \sum_{i=1}^M W_i^{(p-1)} \frac{\eta_p(\boldsymbol{\theta}_i^{(p-1)})}{\eta_{p-1}(\boldsymbol{\theta}_i^{(p-1)})},$$

giving the estimate of the marginal likelihood

$$\widehat{p}(\mathbf{y}) = \prod_{p=1}^P \frac{\widehat{Z}_p}{Z_{p-1}}. \quad (4)$$

## 3 Methodology

### 3.1 Sequence of target densities

Suppose that  $\mathbf{y}_k, k = 1, \dots, n$ , are independent given  $\boldsymbol{\theta}$  so that the likelihood and log-likelihood can be written as

$$L(\boldsymbol{\theta}) = \prod_{k=1}^n p(\mathbf{y}_k|\boldsymbol{\theta}) \quad \text{and} \quad \ell(\boldsymbol{\theta}) = \sum_{k=1}^n \ell_k(\boldsymbol{\theta}), \quad (5)$$

where  $\ell_k(\boldsymbol{\theta}) = \log p(\mathbf{y}_k|\boldsymbol{\theta})$ . We are concerned with the case where the log-likelihood is computationally very costly, because  $n$  is so large that repeatedly computing this sum is impractical, or  $n$  is moderately large but each term is expensive to evaluate.

Quiroz et al. (2019) propose to subsample  $m$  observations and estimate  $L(\boldsymbol{\theta})$  from an unbiased estimator  $\widehat{\ell}_m(\boldsymbol{\theta})$  of  $\ell(\boldsymbol{\theta})$

$$\widehat{L}(\boldsymbol{\theta}) = \exp \left( \widehat{\ell}_m(\boldsymbol{\theta}) - \frac{1}{2} \widehat{\sigma}_m^2(\boldsymbol{\theta}) \right), \quad (6)$$

where  $\hat{\sigma}_m^2(\boldsymbol{\theta})$  is an estimate of  $\sigma^2(\boldsymbol{\theta}) = \mathbb{V}(\hat{\ell}_m(\boldsymbol{\theta}))$ . The motivation for (6) is that  $\exp(\hat{\ell}_m(\boldsymbol{\theta}) - \sigma^2(\boldsymbol{\theta})/2)$  is an unbiased estimator of  $L(\boldsymbol{\theta})$  when  $\hat{\ell}_m(\boldsymbol{\theta})$  is normal (Ceperley and Dewing, 1999). We note that by the central limit theorem,  $\hat{\ell}_m(\boldsymbol{\theta})$  is likely to be normal for moderate  $m$  when  $n$  is large even if  $m$  is a small fraction of  $n$ . More generally, (6) is an unbiased estimator for  $L_{(m,n)}(\boldsymbol{\theta}) := \mathbb{E}(\hat{L}(\boldsymbol{\theta}))$ , which we call the perturbed likelihood. The expectation with respect to the subsampling indices  $\mathbf{u}$  is discussed below. Quiroz et al. (2019) show that when using the control variate in Section 3.2 in the estimator  $\hat{\ell}_m(\boldsymbol{\theta})$ , and under some extra plausible assumptions, the fractional error of the perturbed likelihood is

$$\left| \frac{L_{(m,n)}(\boldsymbol{\theta}) - L(\boldsymbol{\theta})}{L(\boldsymbol{\theta})} \right| = O\left(\frac{1}{nm^2}\right).$$

Our approach is based on extending the target at the  $p$ th density, i.e.  $\pi_p(\boldsymbol{\theta})$  in (3), to include the set of subsampling indices  $\mathbf{u} = (u_1, \dots, u_m)$ , where  $\mathbf{u} \in \mathcal{U} \subset \{1, \dots, n\}^m$  when sampling data observations with replacement. Let  $\hat{L}_p(\boldsymbol{\theta})$  be an estimator of the tempered likelihood  $L(\boldsymbol{\theta})^{a_p}$ . Similarly to Quiroz et al. (2019), we can unbiasedly estimate  $a_p \ell(\boldsymbol{\theta})$  with  $a_p \hat{\ell}(\boldsymbol{\theta})$ , and since  $\mathbb{V}(a_p \hat{\ell}(\boldsymbol{\theta})) = a_p^2 \sigma^2(\boldsymbol{\theta})$  and motivated by (6), we propose the annealed likelihood estimator

$$\hat{L}_p(\boldsymbol{\theta}) = \exp\left(a_p \hat{\ell}_m(\boldsymbol{\theta}) - \frac{1}{2} a_p^2 \hat{\sigma}_m^2(\boldsymbol{\theta})\right). \quad (7)$$

The extended target at the  $p$ th density is

$$\bar{\pi}_p(\boldsymbol{\theta}, \mathbf{u}) \propto \hat{L}_p(\boldsymbol{\theta}) p(\boldsymbol{\theta}) p(\mathbf{u}) = \exp\left(a_p \hat{\ell}_m(\boldsymbol{\theta}) - \frac{1}{2} a_p^2 \hat{\sigma}_m^2(\boldsymbol{\theta})\right) p(\boldsymbol{\theta}) p(\mathbf{u}), \quad (8)$$

where  $p(\mathbf{u})$  is the density of  $\mathbf{u}$  (or, more correctly, a probability mass function since  $\mathbf{u}$  is discrete). At the final annealing step, (8) becomes  $\bar{\pi}_P(\boldsymbol{\theta}, \mathbf{u}) \propto \hat{L}(\boldsymbol{\theta}) p(\boldsymbol{\theta}) p(\mathbf{u})$ , which is the target considered in Quiroz et al. (2019). Quiroz et al. (2019) show that the perturbed marginal density for  $\boldsymbol{\theta}$ ,  $\pi_{(m,n)}(\boldsymbol{\theta}) = \int_{\mathcal{U}} \bar{\pi}_P(\boldsymbol{\theta}, \mathbf{u}) d\mathbf{u}$  converges in the total variation metric to  $\pi(\boldsymbol{\theta})$  at the rate  $O(1/(nm^2))$ . Hence, our proposed approach is approximate but can be very accurate, while also scaling well with respect to the subsample size. For example, if we take  $m = O(\sqrt{n})$ , then by Quiroz et al. (2019, Part (i) of Theorem 1)

$$\int_{\boldsymbol{\Theta}} |\pi_{(m,n)}(\boldsymbol{\theta}) - \pi(\boldsymbol{\theta})| d\boldsymbol{\theta} = O\left(\frac{1}{n^2}\right).$$

Moreover, suppose that  $\varphi(\boldsymbol{\theta})$  is a scalar function with finite second moment. Then,

by Quiroz et al. (2019, Part (ii) of Theorem 1)

$$\left| \mathbb{E}_{\pi_{(m,n)}}(\varphi(\boldsymbol{\theta})) - \mathbb{E}_{\pi}(\varphi(\boldsymbol{\theta})) \right| = O\left(\frac{1}{n^2}\right).$$

Thus, the approximation obtained by our approach converges to the posterior (in total variation norm) at a very fast rate as do the posterior moment estimates. Sections 4 and 5 confirm empirically that we obtain very accurate estimates in most of our applications, even for an  $m$  very small relative to  $n$ .

### 3.2 Efficient estimator of the log-likelihood

Quiroz et al. (2019) propose estimating  $\ell(\boldsymbol{\theta})$  in (5) by the unbiased difference estimator,

$$\widehat{\ell}_m(\boldsymbol{\theta}) = \sum_{k=1}^n q_k(\boldsymbol{\theta}) + \frac{n}{m} \sum_{j=1}^m \ell_{u_j}(\boldsymbol{\theta}) - q_{u_j}(\boldsymbol{\theta}), \quad u_j \in \{1, \dots, n\} \text{ iid}, \quad (9)$$

where

$$\Pr(u_j = k) = \frac{1}{n} \text{ for all } k = 1, \dots, n \text{ and } j = 1, \dots, m,$$

and  $q_k(\boldsymbol{\theta}) \approx \ell_k(\boldsymbol{\theta})$  are control variates. The estimator is based on writing

$$\ell(\boldsymbol{\theta}) = \sum_{k=1}^n q_k(\boldsymbol{\theta}) + \sum_{k=1}^n d_k(\boldsymbol{\theta}) = q(\boldsymbol{\theta}) + d(\boldsymbol{\theta}),$$

with  $d_k(\boldsymbol{\theta}) = \ell_k(\boldsymbol{\theta}) - q_k(\boldsymbol{\theta})$ ,  $q(\boldsymbol{\theta}) = \sum_k q_k(\boldsymbol{\theta})$ , and  $d(\boldsymbol{\theta}) = \sum_k d_k(\boldsymbol{\theta})$ . The last term on the right hand side of (9) is an unbiased estimator of  $d(\boldsymbol{\theta})$ . We now discuss a choice of control variates due to Bardenet et al. (2017), which computes  $q(\boldsymbol{\theta})$  in  $O(1)$  time. Hence, the cost of computing the estimator is  $O(m)$  and we can take  $m = O(\sqrt{n})$  in order to achieve the convergence rates  $O(1/n^2)$  for both the perturbed density and its moments as discussed in Section 3.1.

Let  $\bar{\boldsymbol{\theta}}$  be an estimate of posterior location, for example the posterior mean, obtained from a current particle cloud from  $\bar{\pi}_p(\boldsymbol{\theta}, \mathbf{u})$ . A second order Taylor series expansion of the log-density around  $\bar{\boldsymbol{\theta}}$  is

$$\ell_k(\boldsymbol{\theta}) = \ell_k(\bar{\boldsymbol{\theta}}) + \nabla_{\boldsymbol{\theta}} \ell_k(\bar{\boldsymbol{\theta}})^\top (\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}) + \frac{1}{2} (\boldsymbol{\theta} - \bar{\boldsymbol{\theta}})^\top (\nabla_{\boldsymbol{\theta}\boldsymbol{\theta}^\top}^2 \ell_k(\bar{\boldsymbol{\theta}})) (\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}) + o(\|\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}\|),$$

where  $o(\delta)$  means that  $o(\delta)/\delta \rightarrow 0$  as  $\delta \rightarrow 0$ . We approximate  $\ell_k(\boldsymbol{\theta})$  by

$$q_k(\boldsymbol{\theta}) = \ell_k(\bar{\boldsymbol{\theta}}) + \nabla_{\boldsymbol{\theta}} \ell_k(\bar{\boldsymbol{\theta}})^\top (\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}) + \frac{1}{2} (\boldsymbol{\theta} - \bar{\boldsymbol{\theta}})^\top (\nabla_{\boldsymbol{\theta}\boldsymbol{\theta}^\top}^2 \ell_k(\bar{\boldsymbol{\theta}})) (\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}).$$

Then,

$$q(\boldsymbol{\theta}) = A(\bar{\boldsymbol{\theta}}) + B(\bar{\boldsymbol{\theta}})(\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}) + \frac{1}{2}(\boldsymbol{\theta} - \bar{\boldsymbol{\theta}})^\top C(\bar{\boldsymbol{\theta}})(\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}),$$

where

$$A(\bar{\boldsymbol{\theta}}) = \sum_k \ell_k(\bar{\boldsymbol{\theta}}), B(\bar{\boldsymbol{\theta}}) = \sum_k \nabla_{\boldsymbol{\theta}} \ell_k(\bar{\boldsymbol{\theta}})^\top \text{ and } C(\bar{\boldsymbol{\theta}}) = \sum_k \nabla_{\boldsymbol{\theta}\boldsymbol{\theta}^\top}^2 \ell_k(\bar{\boldsymbol{\theta}}).$$

The sums  $A(\bar{\boldsymbol{\theta}})$ ,  $B(\bar{\boldsymbol{\theta}})$ , and  $C(\bar{\boldsymbol{\theta}})$  are computed only once at every stage of the SMC, regardless of the number of particles. Then, for each particle, estimating  $d(\boldsymbol{\theta})$  by  $\hat{d}_m(\boldsymbol{\theta}) = (n/m) \sum_j d_{u_j}(\boldsymbol{\theta})$  is computed in  $O(m)$  time, and so is (9) because  $q(\boldsymbol{\theta})$  is  $O(1)$ . We estimate  $\sigma^2(\boldsymbol{\theta}) = \mathbb{V}(\hat{\ell}_m(\boldsymbol{\theta}))$  by

$$\hat{\sigma}_m^2(\boldsymbol{\theta}) = \frac{n^2}{m^2} \sum_{j=1}^m (d_{u_j}(\boldsymbol{\theta}) - \bar{d}_{\mathbf{u}}(\boldsymbol{\theta}))^2,$$

where  $\bar{d}_{\mathbf{u}}(\boldsymbol{\theta})$  denotes the mean of the  $d_{u_j}$  for the sample  $\mathbf{u} = (u_1, \dots, u_m)$ . The estimate  $\hat{\sigma}_m^2(\boldsymbol{\theta})$  comes at virtually no cost since it involves terms that are already computed when obtaining  $\hat{d}_m(\boldsymbol{\theta})$ .

### 3.3 The reweighting and resampling steps

The initial particle cloud and weights are now  $\{\boldsymbol{\theta}_{1:M}^{(0)}, \mathbf{u}_{1:M}^{(0)}, W_{1:M}^{(0)}\}$ , obtained by generating the  $\{\boldsymbol{\theta}_{1:M}^{(0)}, \mathbf{u}_{1:M}^{(0)}\}$  from  $p(\boldsymbol{\theta})$  and  $p(\mathbf{u})$ , and assigning equal weights, i.e.,  $W_{1:M}^{(0)} = 1/M$ . The weighted particles  $\{\boldsymbol{\theta}_{1:M}^{(p-1)}, \mathbf{u}_{1:M}^{(p-1)}, W_{1:M}^{(p-1)}\}$  at the  $(p-1)$ st stage are a sample from  $\bar{\pi}_{p-1}(\boldsymbol{\theta}, \mathbf{u})$  and are propagated to  $\bar{\pi}_p(\boldsymbol{\theta}, \mathbf{u})$ , by updating the weights  $W_{1:M}^{(p)} = w_{1:M}^{(p)} / \sum_{i=1}^M w_i^{(p)}$ , where

$$w_i^{(p)} = W_i^{(p-1)} \exp\left(\left(a_p - a_{p-1}\right) \hat{\ell}_m(\boldsymbol{\theta}_i^{(p-1)}) - \frac{1}{2} \left(a_p^2 - a_{p-1}^2\right) \hat{\sigma}_m^2(\boldsymbol{\theta}_i^{(p-1)})\right).$$

The particles  $\{\boldsymbol{\theta}_{1:M}^{(p-1)}, \mathbf{u}_{1:M}^{(p-1)}\}$  are then resampled using the weights  $W_{1:M}^{(p)}$  to obtain the equally-weighted particles  $\{\boldsymbol{\theta}_{1:M}^{(p)}, \mathbf{u}_{1:M}^{(p)}\}$ .

### 3.4 The Markov move step

The Markov move step uses Hamiltonian dynamics to propose distant particle moves and data subsampling to speed up the computation of the dynamics. Similarly to Section 2.1, the Markov move is designed to leave each of the sequence target densities

$\bar{\pi}_p(\boldsymbol{\theta}, \mathbf{u})$ , for  $p = 0, \dots, P$ , invariant. Algorithm 1 describes the Markov move step and is divided into two parts to accommodate subsampling. See Dang et al. (2019) for the details.

---

**Algorithm 1** Single Markov move with a kernel invariant for  $\bar{\pi}_p(\boldsymbol{\theta}, \mathbf{u})$  in (8).

---

For  $i = 1, \dots, M$ ,

1. Sample  $\mathbf{u}_i | \boldsymbol{\theta}_i, \mathbf{y}$ : Propose  $\mathbf{u}_i^* \sim p(\mathbf{u})$ , and set  $\mathbf{u}_i = \mathbf{u}_i^*$ , with probability

$$\alpha_{\mathbf{u}} = \min \left( 1, r := \frac{\exp \left( a_p \widehat{\ell}_m(\boldsymbol{\theta}_i, \mathbf{u}_i^*) - \frac{1}{2} a_p^2 \widehat{\sigma}_m^2(\boldsymbol{\theta}_i, \mathbf{u}_i^*) \right)}{\exp \left( a_p \widehat{\ell}_m(\boldsymbol{\theta}_i, \mathbf{u}_i) - \frac{1}{2} a_p^2 \widehat{\sigma}_m^2(\boldsymbol{\theta}_i, \mathbf{u}_i) \right)} \right), \quad (10)$$

The proposal  $\mathbf{u}_i^*$  is independent of the current value of  $\mathbf{u}_i$ , so the difference between the log of the numerator and log of the denominator of the ratio  $r$  in (10) can be highly variable. This move might get stuck when the denominator is significantly overestimated. A remedy is to induce a high correlation between the log of the estimated annealed likelihood at the current and proposed draws in (10). This can be achieved either through correlating the  $\mathbf{u}$  as in Deligiannidis et al. (2018) (see Quiroz et al. 2019 for discrete  $\mathbf{u}$ ) or by block updates of  $\mathbf{u}$  as in Tran et al. (2017); Quiroz et al. (2018a). We implement the block updates with  $G$  blocks, which gives an approximate correlation  $1 - \frac{1}{G}$ .

2. Sample  $\boldsymbol{\theta}_i | \mathbf{u}_i, \mathbf{y}$ : Given a subset of data  $\mathbf{u}_i$ , we move the particle  $\boldsymbol{\theta}_i$  using a Hamiltonian Monte Carlo (HMC) proposal in a Metropolis-Hastings (MH) algorithm. This becomes a standard HMC move for a given subset  $\mathbf{u}$ .

Note that the above is a Gibbs update of  $\boldsymbol{\theta}_i, \mathbf{u}_i | \mathbf{y}$ . The MH within Gibbs performed in Step 1. is valid (Johnson et al., 2013) and so is the HMC within Gibbs (Neal, 2011) in Step 2. Therefore, this kernel has  $\bar{\pi}_p(\boldsymbol{\theta}, \mathbf{u})$  as its invariant distribution. Dang et al. (2019) previously proposed an MCMC version of this algorithm.

---

Algorithm 2 summarizes our approach. We follow Buchholz et al. (2018) who develop a tuning procedure for the mass matrix, the step size and the number of leapfrog steps within an SMC framework. The number of Markov moves  $R$  is tuned by increasing it until 90% of the product of componentwise autocorrelation of the particles drops below a threshold; see Buchholz et al. (2018) for more details.

---

**Algorithm 2** Subsampling Sequential Monte Carlo
 

---

1. Sample the particles  $\{\boldsymbol{\theta}_i^{(0)}, \mathbf{u}_i^{(0)}\}$  from the prior densities  $p(\boldsymbol{\theta})$  and  $p(\mathbf{u})$  and give all particles equal weights,  $W_i = 1/M$ ,  $i = 1, \dots, M$ .
2. While the tempering sequence  $a_p \neq 1$  do
  - (a) Set  $p \leftarrow p + 1$
  - (b) Find  $a_p$  adaptively to maintain the ESS around  $\text{ESS}_{\text{target}}$  (Section 2.2).
  - (c) Reweighting: compute the unnormalized weights

$$\begin{aligned}
 w_i^{(p)} &= W_i^{(p-1)} \frac{\eta_{a_p}(\boldsymbol{\theta}_i^{(p-1)}, \mathbf{u}_i^{(p-1)})}{\eta_{a_{p-1}}(\boldsymbol{\theta}_i^{(p-1)}, \mathbf{u}_i^{(p-1)})} \\
 &= W_i^{(p-1)} \exp\left((a_p - a_{p-1}) \widehat{\ell}_m(\boldsymbol{\theta}_i^{(p-1)}) - \frac{1}{2} (a_p^2 - a_{p-1}^2) \widehat{\sigma}_m^2(\boldsymbol{\theta}_i^{(p-1)})\right),
 \end{aligned}$$

and normalize as  $W_i^{(p)} = w_i / \sum_{i'=1}^M w_{i'}$ ,  $i = 1, \dots, M$ .

- (d) Compute  $\bar{\boldsymbol{\theta}}$  as  $\bar{\boldsymbol{\theta}} = \sum_{i=1}^M W_i^{(p)} \boldsymbol{\theta}_i^{(p-1)}$  and then obtain

$$\sum_{k=1}^n \ell_k(\bar{\boldsymbol{\theta}}), \quad \sum_{k=1}^n \nabla_{\boldsymbol{\theta}} \ell_k(\bar{\boldsymbol{\theta}}), \quad \sum_{k=1}^n \nabla_{\boldsymbol{\theta}\boldsymbol{\theta}^\top}^2 \ell_k(\bar{\boldsymbol{\theta}})$$

and the mass matrix  $\mathbf{H} = \Sigma^{-1}(\bar{\boldsymbol{\theta}})$ , where  $\Sigma$  is the sample covariance matrix of current particles.

- (e) Resample the particles  $\{\boldsymbol{\theta}_i^{(p-1)}, \mathbf{u}_i^{(p-1)}\}_{i=1}^M$  using the weights  $\{W_i^{(p)}\}_{i=1}^M$  to obtain resampled particles  $\{\boldsymbol{\theta}_i^{(p)}, \mathbf{u}_i^{(p)}\}_{i=1}^M$  and set  $W_i^{(p)} = 1/M$ .
  - (f) Apply  $R$  Markov moves to each particle  $\boldsymbol{\theta}_i^{(p)}, \mathbf{u}_i^{(p)}$  using Algorithm 1.
- 

### 3.5 Marginal likelihood estimation

Our approach naturally extends that of Section 2.3 by considering the augmented target density  $\bar{\pi}_p(\boldsymbol{\theta}, \mathbf{u})$  in (8). Define

$$\gamma_p(\boldsymbol{\theta}, \mathbf{u}) = \frac{\eta_p(\boldsymbol{\theta}, \mathbf{u})}{\eta_{p-1}(\boldsymbol{\theta}, \mathbf{u})}.$$

Then

$$\begin{aligned} \int_{\mathcal{U}} \int_{\Theta} \gamma_p(\boldsymbol{\theta}, \mathbf{u}) \pi_{p-1}(\boldsymbol{\theta}, \mathbf{u}) d\boldsymbol{\theta} d\mathbf{u} &= \int_{\mathcal{U}} \int_{\Theta} \frac{\eta_p(\boldsymbol{\theta}, \mathbf{u})}{\eta_{p-1}(\boldsymbol{\theta}, \mathbf{u})} \frac{\eta_{p-1}(\boldsymbol{\theta}, \mathbf{u})}{Z_{p-1}} p(\boldsymbol{\theta}) p(\mathbf{u}) d\boldsymbol{\theta} d\mathbf{u} \\ &= \frac{Z_p}{Z_{p-1}}. \end{aligned}$$

Thus, if  $\{\boldsymbol{\theta}_{1:M}^{(p-1)}, \mathbf{u}_{1:M}^{(p-1)}, W_{1:M}^{(p-1)}\}$  at the  $(p-1)$ st sequence is an approximate sample from  $\bar{\pi}_{a_{p-1}}(\boldsymbol{\theta}, \mathbf{u})$ , we estimate the ratio  $Z_p/Z_{p-1}$  by

$$\frac{\widehat{Z}_p}{Z_{p-1}} = \sum_{i=1}^M W_i^{(p-1)} \frac{\eta_p(\boldsymbol{\theta}_i^{(p-1)}, \mathbf{u}_i^{(p-1)})}{\eta_{p-1}(\boldsymbol{\theta}_i^{(p-1)}, \mathbf{u}_i^{(p-1)})},$$

and the marginal likelihood estimate is obtained using this expression in (4).

### 3.6 Efficient memory management by data subsampling

We now explain in detail how data subsampling helps to parallelize the computing in terms of efficient memory utilization. Suppose first that we perform standard SMC (using all the data) and that we parallelise using  $N$  workers, so that each worker deals, on average, with  $M/N$  particles. Then, for each stage  $p$ , the computations performed for each particle require repeated likelihood evaluations (using all  $n$  data) when applying  $R$  Markov move steps. Hence, each worker needs to have access to the full dataset.

Suppose now that we use our data subsampling approach in the same setting using  $M/N$  particles for each of the  $N$  workers. Then, at the beginning of each stage  $p$  of the algorithm, we still require a full data evaluation for computing  $A(\bar{\boldsymbol{\theta}})$ ,  $B(\bar{\boldsymbol{\theta}})$  and  $C(\bar{\boldsymbol{\theta}})$  in Section 3.2. However, at each  $p$ , we can now subsample the data according to  $\mathbf{u}_i^{(p)}$  for each particle and subsequently perform the  $R$  Markov move steps, which now require repeated evaluations of the estimated annealed likelihood (using  $m \ll n$  observations) and in addition  $A(\bar{\boldsymbol{\theta}})$ ,  $B(\bar{\boldsymbol{\theta}})$  and  $C(\bar{\boldsymbol{\theta}})$ . Now each worker needs to have access only to the subsampled dataset, as well as  $A(\bar{\boldsymbol{\theta}})$ ,  $B(\bar{\boldsymbol{\theta}})$  and  $C(\bar{\boldsymbol{\theta}})$ . However, these are only summaries of the full dataset and are therefore very memory efficient.

We are aware that parallelization of SMC methods is not straightforward to do efficiently when resampling occurs often (Murray et al., 2016; Lee et al., 2010). We note that in all our applications the number of annealing steps is relative small and therefore resampling does not really affect the efficiency of our algorithm. In applications where resampling occurs more frequently, both SMC methods can benefit from the ideas in Heine et al. (2019) and Guldás et al. (2015). Moreover, the reweighting

and the computationally expensive Markov move steps of our algorithm are easily parallelised for each SMC sample because the computations required for each sample are independent of those of the other samples. Subsampling therefore does not affect the parallelisation of the algorithm because only the part of the data specified by the particles  $u_i$  are sent to each worker and the  $u_i$  are independent of each other.

## 4 Evaluations

### 4.1 Experiments

We now evaluate the methodology through the following experiments.

- *Experiment 1: Evaluating the usefulness of the Hamiltonian Monte Carlo kernel.*

We show the effectiveness of a HMC kernel for the Markov move step compared to random walk and MALA kernels.

- *Experiment 2: Evaluating the speed and the accuracy of the marginal likelihood and the approximate posterior density when the posterior is unimodal.*

We show that the subsampling approach is accurate by comparing the estimates of the marginal likelihood and posterior densities to those obtained by the full data SMC (representing the gold standard).

- *Experiment 3: Evaluating the speed and the accuracy of the marginal likelihood and the approximate posterior density when the posterior is non-Gaussian.*

We use the subsampling approach when the posterior is bimodal or skewed and show that the method still performs well.

- *Experiment 4: Evaluating the effect of the accuracy of the control variate.*

We show that the subsampling approach can be made faster by using a first order control variate instead of the second order alternative. This experiment also shows the effect of inaccurate likelihood estimates on the performance of our method.

All the SMC algorithms are tuned as in Buchholz et al. (2018) using 280 particles, a choice motivated by our cluster with 28 cores with each core dealing (on average) with 10 particles. The only exception is the first scenario in Experiment 3 where we use 420 particles for both algorithms to better capture the multimodal posterior. We repeated each experiment 10 times to compute the standard error of the log marginal likelihood estimator. Experiments 1, 2 and the bankruptcy application in

Section 5 were done using the Australia NCI High Performance Computing System Raijin<sup>1</sup>. Experiments 3 and 4 were done using the University of New South Wales computational cluster Katana<sup>2</sup>.

We remark that the choice of priors can affect the computational efficiency of SMC methods. In general, a prior that resembles the likelihood requires less tempering steps. However, this is unlikely to influence the comparison between Subsampling SMC and SMC, which is our primary concern.

## 4.2 Experiment 1: Evaluating the Markov move kernel

We first consider a logistic regression to evaluate how effectively the Hamiltonian Monte Carlo Markov move step is compared to the random walk and MALA kernels. The model for the response  $y_i \in \{0, 1\}$  given a  $d \times 1$  set of covariates and parameters is

$$p(y_i | \mathbf{x}_i, \boldsymbol{\theta}) = \frac{\exp(y_i \mathbf{x}_i^\top \boldsymbol{\theta})}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\theta})}.$$

We fit this model to the HIGGS dataset (Baldi et al., 2014), having  $n = 11,000,000$  observations and 28 covariates. The response is “detected particle” and 21 of the covariates are kinematic properties measured by particle detectors, while 7 are high-level features to capture non-linearities. This means that  $d = 29$ , including the intercept. We take the prior  $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , where  $\mathbf{I}_d$  is the  $d \times d$  identity matrix and follow Buchholz et al. (2018) to set the tuning parameters, including the number of Markov moves  $R$ . The mass matrix in both HMC and MALA is  $\widehat{\Sigma}^{-1}$ , which is the estimated inverse covariance matrix of the tempered posterior. We note that each step in the sequence has a corresponding estimate of this inverse covariance matrix, obtained using the corresponding particles from that step. For the random walk, the optimal scaling  $(2.38^2/d)\widehat{\Sigma}$  (Roberts et al., 1997) resulted in numerical errors, so that we decreased it by a factor of 10.

Table 1 summarizes the results obtained using the second order control variate in Section 3.2. The log-likelihood estimator has  $m = 5,000$  subsamples and the block-pseudo marginal is carried out using  $G = 100$ . Clearly, the Hamiltonian approach is computationally faster because it needs to take a smaller number of Markov steps  $R$ . The table also shows that the log of the estimate of marginal likelihood is very similar for all methods. The rest of the article uses the HMC kernel.

---

<sup>1</sup><https://nci.org.au/our-systems/hpc-systems>

<sup>2</sup><https://research.unsw.edu.au/katana>

Table 1: Comparing the performances of three kernels for the Markov move, Hamiltonian Monte Carlo (HMC), Metropolis Adjusted Langevin Algorithm (MALA) and Random Walk (RW). The table shows the log of the estimate of the marginal likelihood (with standard error in brackets), the CPU time, the number of annealing steps  $P$  (tuned to maintain  $ESS \approx 0.8M$ ) and the number of Markov moves  $R$  (tuned as in Buchholz et al., 2018). The results are for the logistic regression model estimated using the HIGGS data and  $M = 280$  particles. All methods use the second order control variate in Section 3.2. The results are averaged over 10 runs, which are used to compute the standard error of the estimator.

	log marginal likelihood	CPU time (hrs)	$P$	$R$
HMC	-7,013,460.90 (0.32)	2.31	106	5
MALA	-7,013,462.49 (0.26)	4.77	106	20
RW	-7,013,461.43 (0.32)	33.43	106	200

### 4.3 Experiment 2: Evaluating speed and accuracy of Subsampling SMC on unimodal targets

This section compares Subsampling SMC with full data SMC. Such a comparison is infeasible for the full HIGGS dataset because it is too large; the full dataset needs to be available at each worker (we use 28) as explained in Section 3.6, in order to compute the likelihood together with its gradient and Hessian, which would quickly consume the RAM of the computer. Instead, we consider the following two models.

**Student-t regression.** We consider a univariate Student-t regression

$$y_i = \mathbf{x}_i^\top \boldsymbol{\theta} + e_i, e_i \sim t_5,$$

where  $t_5$  is the Student-t distribution with 5 degrees of freedom. We generated a simulated dataset with  $n = 500,000$  observations and  $d = 50$  covariates. The covariates were generated so that their marginal variances are 1 and their pairwise correlations are 0.9. The parameters  $\boldsymbol{\theta}$  were simulated independently from a Uniform( $-5, 5$ ) distribution; the prior for  $\boldsymbol{\theta}$  is  $\mathcal{N}(\mathbf{0}, 10\mathbf{I}_d)$ .

**Poisson regression.** We also considered a Poisson regression, where the univariate  $y$  follows a Poisson distribution with an expectation that is log-linear, i.e.

$$y_i | \mathbf{x}_i \sim \text{Poisson}(\exp(\mathbf{x}_i^\top \boldsymbol{\theta})).$$

Table 2: Comparing the performances of Subsampling SMC and full data SMC. The table shows the log of the estimate of the marginal likelihood (with standard error in brackets), the CPU time, the number of annealing steps  $P$  (tuned to maintain  $ESS \approx 0.8M$ ) and the number of Markov moves  $R$  (tuned as in Buchholz et al., 2018). The results are for the Student-t regression and Poisson regression models estimated using the simulated datasets. We use  $M = 280$  particles. All methods use the second order control variate in Section 3.2. The results are averaged over 10 runs, which are used to compute the standard error of the estimator.

	log marginal likelihood	CPU time (hrs)	$P$	$R$
<b><u>Student-t regression</u></b>				
$(n = 500,000, m = 1,200)$				
Full data SMC	-815,775.82 (0.39)	5.92	126	4
Subsampling SMC	-815,773.49 (0.59)	0.57	127	4
<b><u>Poisson regression</u></b>				
$(n = 200,000, m = 500)$				
SMC	-260,888.69 (1.40)	0.94	80	4
Subsampling SMC	-260,887.87 (0.27)	0.14	80	5

We generated  $n = 200,000$  observations with  $d = 30$  covariates, 29 of them simulated from  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{29})$  and the last one is 1. The parameters are simulated independently from  $\text{Uniform}(-0.2, 0.2)$  and are assigned the prior  $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, 0.1\mathbf{I}_d)$ . We found that both Subsampling SMC and SMC were particularly sensitive to the prior choice for the Poisson regression, resulting in numerical overflow for priors that were too diffuse.

For both examples, we used  $G = 100$  blocks and the second order Taylor series control variates and set  $m$  to correspond to a sample fraction of about 0.0025. Table 2 summarizes the results and shows that the subsampling approach is about 6.5 to 10.5 times faster and, moreover, confirms the accuracy of the marginal likelihood estimate of our method. Finally, Figures 1 and 2 show that the marginal posterior densities are very well approximated for both the Student-t regression and the Poisson regression; the same accuracy was obtained for all parameters (not shown).

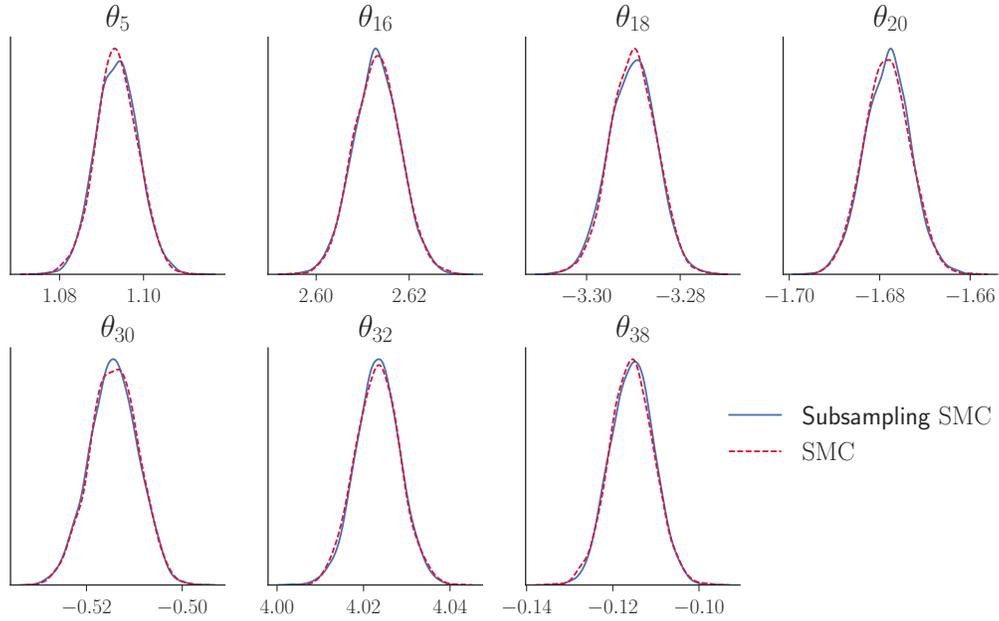


Figure 1: Kernel density estimates of a subset of the marginal posterior densities of  $\theta$  for the Student-t regression model with simulated data. The density estimates are obtained by full data SMC and Subsampling SMC.

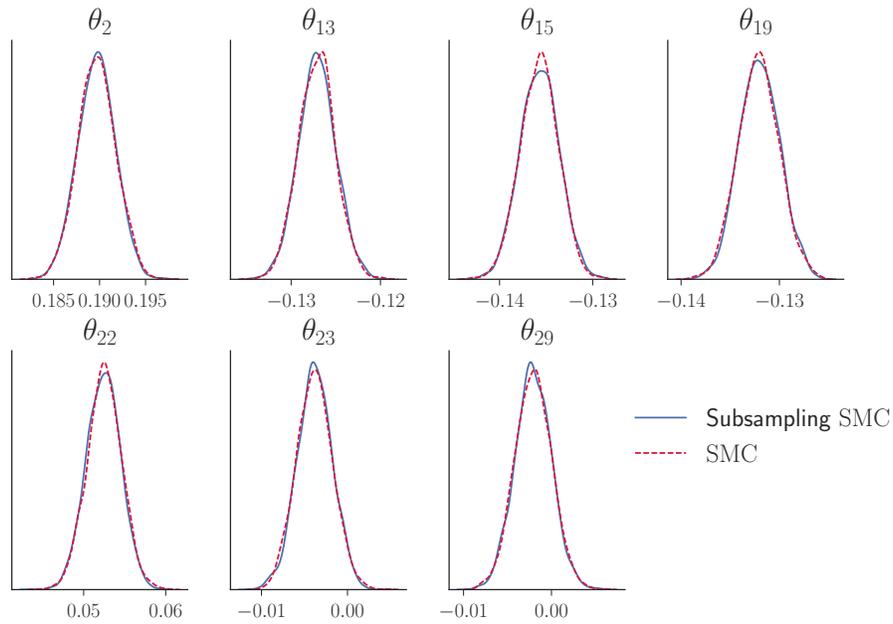


Figure 2: Kernel density estimates of a subset of the marginal posterior densities of  $\theta$  for the Poisson regression model with simulated data. The density estimates are obtained by full data SMC and Subsampling SMC.

## 4.4 Experiment 3: Evaluating speed and accuracy of Subsampling SMC on non-Gaussian targets

To evaluate the performance of Subsampling SMC when the posterior is non-Gaussian, we consider the fixed effects model

$$y_{ij} = \alpha_i + \mathbf{x}_{ij}^\top \boldsymbol{\beta} + e_{ij}, \quad i = 1, \dots, n, \quad j = 1, \dots, n_i, \quad e_{ij} \sim \mathcal{N}(0, \sigma_i^2).$$

For simplicity, we set  $\sigma_i^2 = 1$  for all 10 individuals; two different scenarios are used for the individual fixed effects  $\alpha_i$ ; a) a mixture of normals prior for the  $\alpha_i$ ; b) a truncated normal prior for the  $\alpha_i$ . For each scenario, we generated a dataset of  $n = 10$  individuals, with  $n_1 = \dots = n_5 = 20$  observations and  $n_6 = \dots = n_{10} = 50,000$  observations. The covariates were generated independently from  $\mathcal{N}(0, 1)$ ; the  $\boldsymbol{\beta}$  parameters were generated from  $\mathcal{N}(0, 2^2 \mathbf{I}_{10})$ . The prior for  $\boldsymbol{\beta}$  in both scenarios is  $\mathcal{N}(0, 3^2 \mathbf{I}_{10})$ .

**Mixture of normals prior** The first prior is motivated by a variable selection scenario, where some coefficients may be 0 or very close to 0 and we would like the posterior to set these close to zero. In this experiment, the first 5 individual fixed effects  $\alpha_i$  were generated from  $\mathcal{N}(0.5, 0.05^2)$  and the rest from  $\mathcal{N}(0.5, 0.2^2)$ . For each of the fixed effects  $\alpha_i$  we used a mixture of normals prior

$$p(\alpha_i | w, \sigma_1, \sigma_2) = w \phi(\alpha_i | \sigma_1^2) + (1 - w) \phi(\alpha_i | \sigma_2^2), \quad i = 1, \dots, n;$$

$\phi(\cdot | \sigma^2)$  is the density of the normal distribution with mean 0 and variance  $\sigma^2$ , and we set  $w = 0.8$ ,  $\sigma_1 = 0.1$  and  $\sigma_2 = 3.5$ .

Even though the likelihood for each individual is likely to be unimodal, the prior leads to more complicated posteriors for those individual effects that correspond to subjects with a small number of observations. The likelihood is

$$L(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \prod_{i=1}^n p(\mathbf{y}_i | \alpha_i, \boldsymbol{\beta}), \quad \text{where} \quad p(\mathbf{y}_i | \alpha_i, \boldsymbol{\beta}) = \prod_{j=1}^{n_i} p(y_{ij} | \alpha_i, \boldsymbol{\beta}) \quad (11)$$

is the likelihood for subject  $i$ . The likelihood  $L(\boldsymbol{\alpha}, \boldsymbol{\beta})$  and the annealed likelihood are estimated by estimating the individual likelihoods with subsampling. The subsample size is  $m = 5$ , with no blocking for the first 5 individuals and  $m = 100$  with  $G = 100$  blocks for the remaining 5 individuals. In practice, it is unnecessary to estimate the likelihood for the individuals with few observations since it is relatively cheap computationally to evaluate their full likelihood; however, we do so in our experiment

to gain more knowledge about the effect of subsampling.

We ran Subsampling SMC with second order Taylor series expansions in both scenarios. Table 3 summarizes the results and shows that Subsampling SMC produces similar results to full data SMC but is about 9 times faster. All SMC methods require the maximum number of Markov moves at most temperatures, indicating that the posterior is challenging to explore. Figure 3 shows that even when some of the marginal posteriors ( $\alpha_1, \alpha_3$  and  $\alpha_5$ ) are bimodal, Subsampling SMC is able to capture that and gives the same approximation as full data SMC. As a comparison, we also include in the figure the result from running 10,000 post burn-in iterations of Subsampling MCMC. It is well known that conventional MCMC methods may not be able to sample efficiently from multimodal targets, and in this experiment Subsampling MCMC can detect the posterior modes, but there is still some visible discrepancy between its result and that of full data SMC. We do not show the marginal posterior densities of  $\beta$  which appear to be Gaussian, but confirm that both methods give similar results.

Our method works in this example because the bimodality is caused by the prior and not the likelihood; if the bimodality was caused by the likelihood, different control variates would be necessary since our control variates assume the log-density is quadratic  $\theta$ . We leave the development of more flexible control variates for Subsampling SMC for future research.

**Truncated normal prior** The second scenario is motivated by situations in which there is strong prior knowledge that the coefficients are positive. To create such a situation, the fixed effects  $\alpha_i$  were generated from a truncated normal distribution  $\mathcal{TN}(0.1, 0.1^2)$ . We assigned the prior  $\alpha_i \sim \mathcal{TN}(0, 3^2), i = 1, \dots, 10$  to the individual fixed effects to reflect this prior knowledge. The subsample size is  $m = 20$  (all observations) with no blocking for the first 5 individuals and  $m = 200$  with  $G = 100$  blocks for the 6<sup>th</sup> individual. The remaining 4 individuals have  $m = 100$  with  $G = 100$ . Note that the subsample size affects the variance of the log-likelihood estimator and hence the accuracy of our method; see Quiroz et al. (2019) and Dang et al. (2019) for further discussion. Section 4.5 discusses the results when a smaller subsample size is used for this model.

Table 3 and Figure 4 summarize the results of full data SMC and subsampling SMC. For this example, the truncated normal prior makes the posterior of  $\alpha_1, \dots, \alpha_5$  skewed. These are the effects corresponding to the subjects with few observations. Subsampling SMC seems to experience some difficulties with this challenging target, which is shown by the slightly higher  $P$  and  $R$  values compared to full data SMC. However our method is still slightly faster and produces posterior estimates similar

Table 3: Comparing the performances of Subsampling SMC and full data SMC. The table shows the log of the estimate of the marginal likelihood (with standard error in brackets), the CPU time, the number of annealing steps  $P$  (tuned to maintain  $ESS \approx 0.8M$ ) and the number of Markov moves  $R$  (tuned as in Buchholz et al., 2018, and the maximum number of Markov moves at each temperature is set to be 100). The results are for the fixed effects model estimated using the simulated datasets. All methods use the second order control variate in Section 3.2. The results are averaged over 10 runs, which are used to compute the standard error of the estimator.

	log marginal likelihood	CPU time (hrs)	$P$	$R$
<b>Mixture of normals priors</b>				
$(M = 420)$				
Full data SMC	-354,914.89 (0.78)	14.36	81	99
Subsampling SMC	-354,915.22 (1.18)	1.60	81	100
<b>Truncated normal priors</b>				
$(M = 280)$				
Full data SMC	-354,445.12 (0.26)	0.74	79	5
Subsampling SMC	-354,444.04 (2.2)	0.68	88	20

to the full data SMC (see Figure 4). We do not show the marginal posterior densities of  $\beta$  which appear to be Gaussian, but both methods gave similar results.

#### 4.5 Experiment 4: Evaluating the effect of the control variate

The results above show that the subsampling approach accurately estimates the marginal likelihood and marginal posterior densities using a second order Taylor series expansion. We now study robustness of the results to the quality of the control variates. The first study uses first order Taylor expansions for the control variates for subsampling applied to the logistic regression for the HIGGS dataset in Section 4.2. Table 4 summarizes the results, and confirms that the marginal likelihood estimator remains accurate, and is five times faster than using the second order control variates. Figure 5 shows that the marginal posterior densities remain accurate, we have confirmed similar accuracy for all the parameters.

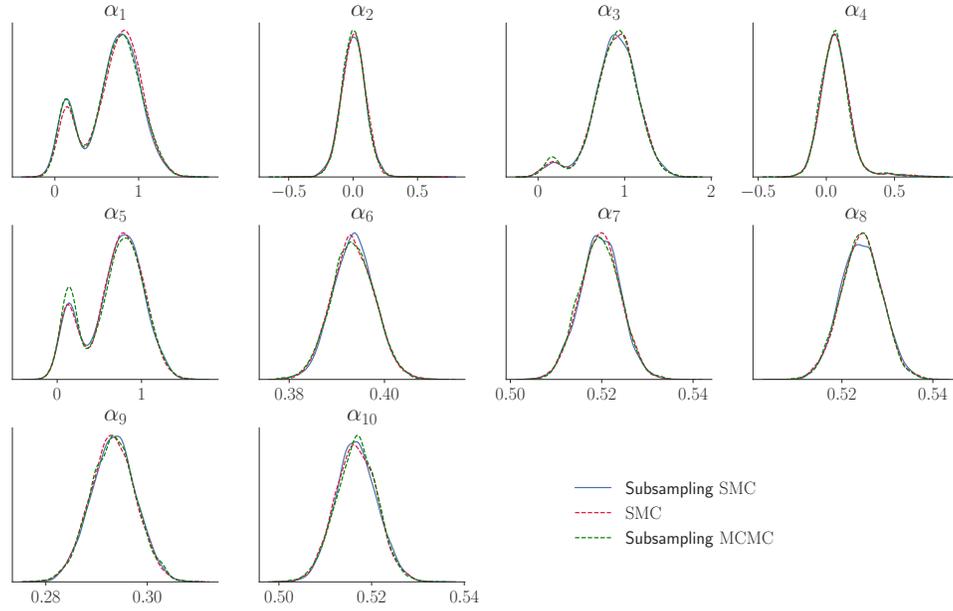


Figure 3: Kernel density estimates of a subset of the marginal posterior densities of  $\alpha$  for the fixed effects model with mixture of normals priors, using simulated data. The density estimates are obtained by full data SMC, Subsampling SMC and Subsampling MCMC.

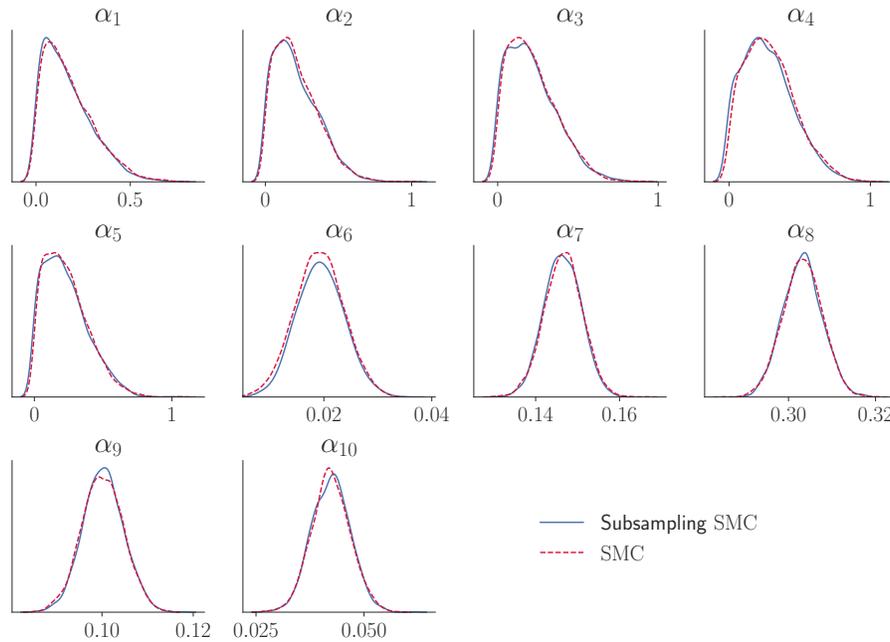


Figure 4: Kernel density estimates of a subset of the marginal posterior densities of  $\alpha$  for the fixed effects model with truncated normal priors, using simulated data. The density estimates are obtained by full data SMC and Subsampling SMC.

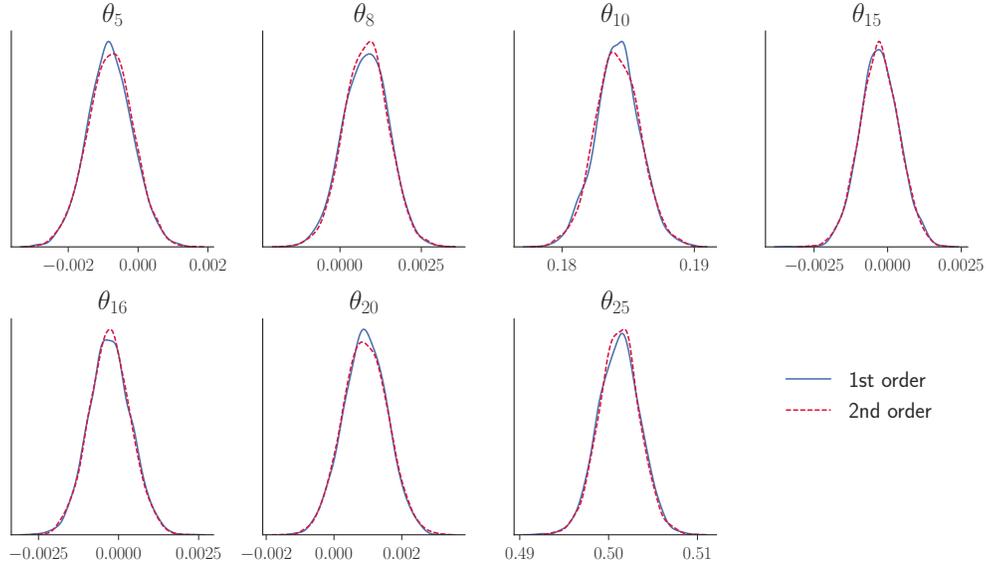


Figure 5: Kernel density estimates of a subset of the marginal posterior densities of  $\theta$  for the logistic model with the HIGGS data. The density estimates are both obtained by Subsampling SMC, using different control variates based on a 1st and 2nd order Taylor series expansion as explained in Section 3.2.

We now present an example where inaccurate likelihood estimates lead to a biased result. We consider again the fixed effects model described in Section 4.4 with the individual effects  $\alpha_i$  having a truncated normal prior,  $p(\alpha_i) \sim \mathcal{TN}(0, 3^2)$ ,  $m = 5$  is used for the first 5 individuals and  $m = 100$  with  $G = 100$  for the remaining 5 individuals.

Table 4 and Figure 6 summarize the results; they show that Subsampling SMC has difficulties exploring the skewed posteriors and gives inaccurate results when  $m$  is too small. Our approach works poorly here because the posteriors for the first 5 individual effects are highly skewed, and their skewness is caused by the truncated prior and the small number of observations. This causes the posterior to be concentrated at the tail of the log-density, where the control variates using a quadratic approximation are inaccurate. Therefore updating  $\bar{\theta}$  by the posterior mean as specified in Algorithm 2 does not produce good control variates, even though the log-density is well-behaved. Our likelihood estimate is inaccurate with high variance even when we use a slightly smaller subsample size compared to the previous section for estimating these skewed posteriors. We leave the development of more flexible control variates and the guidelines to choose an optimal subsample size, especially for complex posteriors, for future research. Finally, Subsampling SMC is not faster than full data SMC here because of the much larger  $P$  and  $R$  chosen by using the adaptive tuning method by Buchholz et al. (2018).

Table 4: Comparing the performance of the less accurate control variate (1st order) to the more accurate control variate (2nd order). The table shows the log of the estimate of the marginal likelihood (with standard errors in brackets), the CPU time, the number of annealing steps  $P$  (tuned to maintain  $ESS \approx 0.8M$ ) and the number of Markov moves  $R$  (tuned as in Buchholz et al. (2018)). The results are for the logistic regression model, estimated with the HIGGS dataset, using  $M = 280$  particles. The results are averaged over 10 runs, which are used to compute the standard error of the estimator.

	log marginal likelihood	CPU time (hrs)	$P$	$R$
<b><u>Logistic regression</u></b>				
1st order	-7,013,461.07 (0.46)	0.47	106	5
2nd order	-7,013,460.90 (0.32)	2.31	106	5
<b><u>Truncated normal priors</u></b>				
Full data SMC	-354,445.12 (0.26)	0.74	79	5
Subsampling SMC	-354,437.34 (4.64)	1.13	141	36

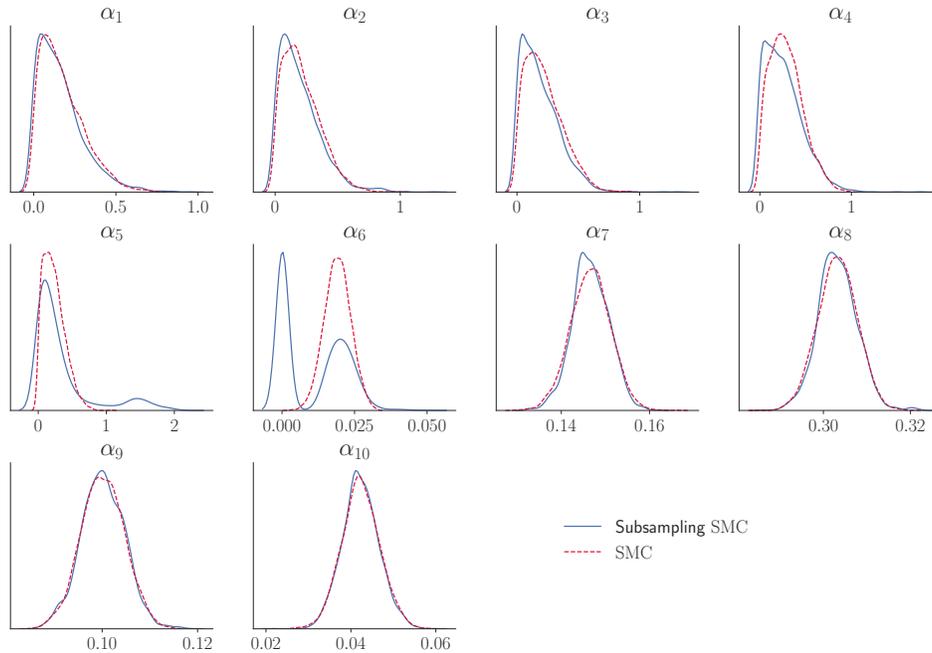


Figure 6: Kernel density estimates of a subset of the marginal posterior densities of  $\alpha$  for the fixed effects model with truncated normal priors, using simulated data. The density estimates are obtained by full data SMC and Subsampling SMC. Subsampling MCMC fails to work on this example and hence its result is not included here.

## 5 Application: Modeling firm bankruptcy nonlinearly

The application of our method for model selection is now illustrated using a Swedish firm bankruptcy dataset containing  $n = 4,748,089$  observations; the response variable is firm default and there are eight firm-specific and macroeconomic covariates, giving 9 covariates, including an intercept. The data is treated as cross-sectional data and the bank status is modeled by the logistic regression discussed in Section 4.2. A generalized additive model is also fitted to the data and is compared to a linear model; a similar prior  $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, 10^2 \mathbf{I}_d)$  is used in both models. We compare the marginal posterior density estimates of Subsampling SMC against those of Subsampling MCMC (Quiroz et al., 2019) as implemented by Dang et al. (2019) and find them nearly indistinguishable. We also compare both methods to the full data MCMC as in Dang et al. (2019). However, it is unclear how to use Subsampling MCMC for model selection. Frequently used methods such as Chib and Jeliazkov (2001) are not useful for Subsampling MCMC since the (perturbed) likelihood cannot be evaluated; this is a major advantage of Subsampling SMC compared to Subsampling MCMC.

We select between model  $\mathcal{M}_1$  which is linear in the data on the logit scale and has 9 coefficients, and model  $\mathcal{M}_2$  which is a semi-parametric additive model on the logit scale and uses B-splines as in Dang et al. (2019); model  $\mathcal{M}_2$  is nonlinear in the data and has 81 coefficients. Non-linear bankruptcy models for this dataset have previously been analyzed in Quiroz and Villani (2013) and Giordani et al. (2014). Given the marginal likelihood estimates, the estimated Bayes Factor (BF) for the non-linear model  $\mathcal{M}_2$  vs the linear model  $\mathcal{M}_1$  is

$$\widehat{\text{BF}}_{21} = \frac{\widehat{\text{Pr}}(\mathbf{y}|\mathcal{M}_2)}{\widehat{\text{Pr}}(\mathbf{y}|\mathcal{M}_1)}; \quad (12)$$

this is also the estimated ratio of posterior model probabilities when the prior model probabilities are equal. We use the strength of evidence guidelines in Jeffreys (1961, p. 438) to choose between the models; Jeffreys considers  $10^{3/2} < \text{BF}_{21} < 10^2$  as very strong evidence for model  $\mathcal{M}_2$  and  $\text{BF}_{21} > 10^2$  as decisive evidence.

The number of blocks was set to  $G = 100$  with the subsample size set to  $m = 3,000$ ; for Subsampling MCMC these tuning parameters were set as in Dang et al. (2019). The estimates from the full data MCMC are considered as the “gold standard” when assessing the accuracy of the algorithms. This was achieved through an MCMC chain of 2,000 post burnin MCMC samples, with the burnin = 1,000 iterations. The MCMC mixed well and we believe that the iterates represent the

Table 5: Log of the estimates of the marginal likelihoods and Bayes factors  $\text{BF}_{21}$  in (12) for selecting between  $\mathcal{M}_1$  and  $\mathcal{M}_2$ . The estimates of the Standard Errors (SE) are in brackets. The SE is computed using the 10 independent parallel runs. The prior probabilities are  $\Pr(\mathcal{M}_1) = \Pr(\mathcal{M}_2) = 1/2$ .

	$\log \hat{p}(\mathbf{y} \mathcal{M}_1)$	$\log \hat{p}(\mathbf{y} \mathcal{M}_2)$	$\hat{\text{BF}}_{21}$
Bankruptcy	-208,517.79 (0.21)	-200,215.13 (6.57)	$\exp(8,302.66)$

posterior adequately .

Table 5 reports the estimated log of the marginal likelihood for both models and the corresponding Bayes factors obtained by Subsampling SMC. The table shows decisively that the non-linear model is superior. We again stress that producing marginal likelihood estimates is very convenient by SMC, whereas it is currently not possible with Subsampling MCMC.

Figure 7 shows the kernel density estimates of the marginal posterior of selected parameters of the non-linear model for the bankruptcy dataset. It is evident that both Subsampling SMC and Subsampling MCMC are very accurate and we have confirmed the accuracy of the kernel density estimates for all the parameters, which we do not show here to save space. Instead, Figure 8 shows the estimated marginal posterior expectations and posterior variances by the two algorithms for all the parameters in the non-linear models. This confirms the accuracy of the estimates of each parameter. We have also confirmed that the kernel density estimates and the estimated marginal posterior expectations and posterior variances are accurate for the linear model (not shown here).

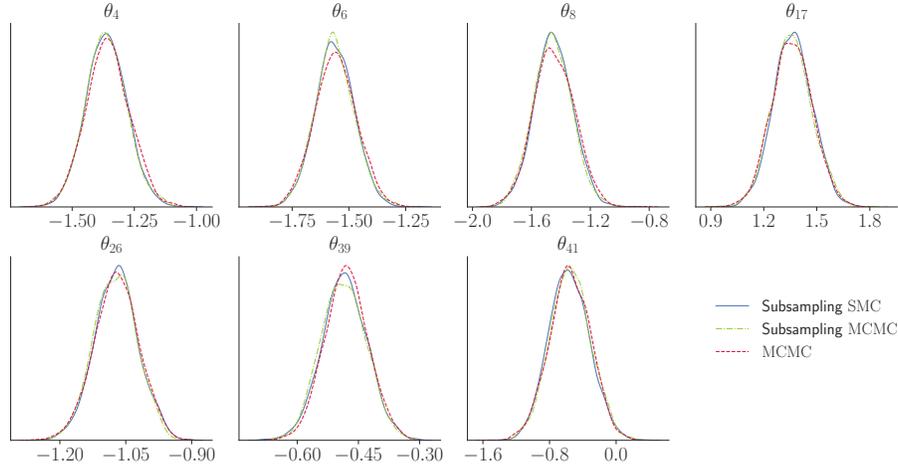


Figure 7: Kernel density estimates of a subset of the marginal posterior densities of  $\theta$  for the logistic model  $\mathcal{M}_2$  for the bankruptcy dataset. The density estimates are obtained by MCMC, Subsampling MCMC and Subsampling SMC. MCMC represents the ground truth.

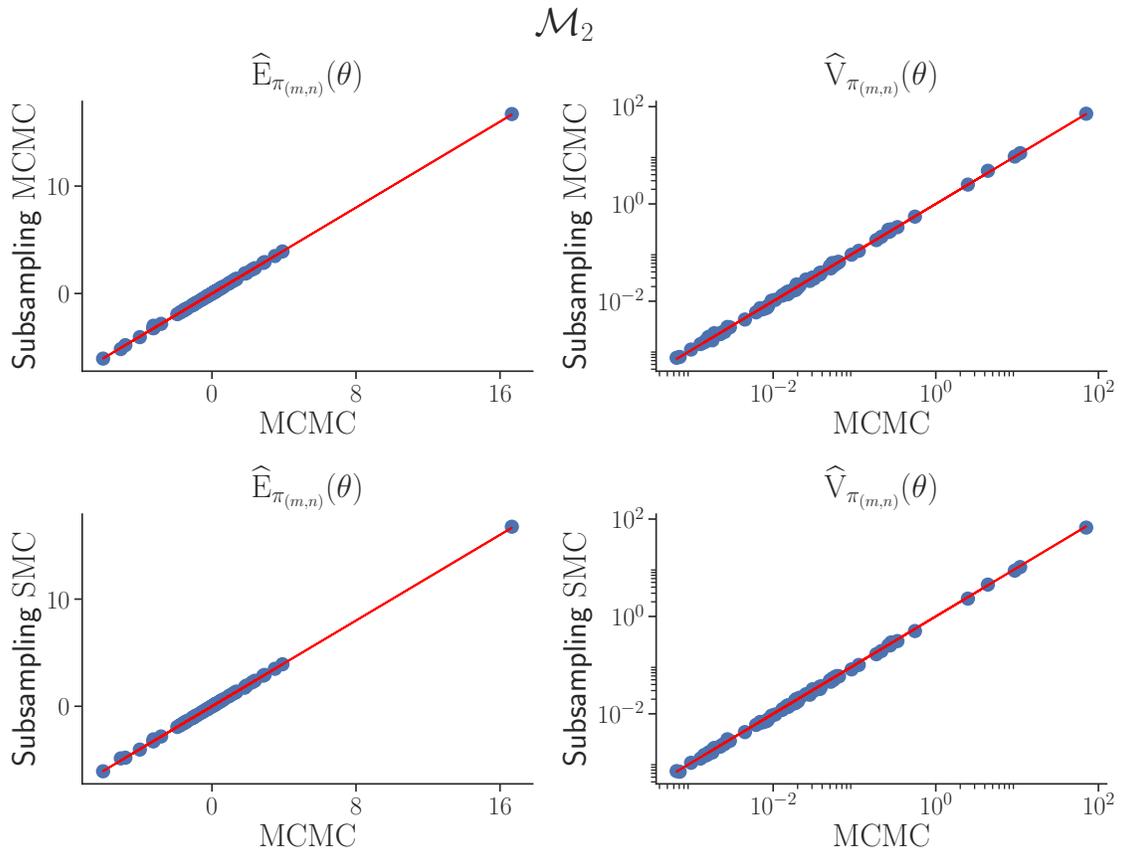


Figure 8: Estimates of marginal posterior means (left panel) and posterior variances (right panel) of  $\theta$  for the logistic model  $\mathcal{M}_2$  for the bankruptcy dataset. The estimates are obtained by Subsampling MCMC and Subsampling SMC and plotted as dots, together with a 45 degree line which corresponds to estimates that are in perfect agreement.

Figure 9 shows that the relationship between the probability of bankruptcy and the covariate Size is not a logistic function (inverse-logit) of the covariate and that the nonlinear model fits the data much better than the linear logistic model.

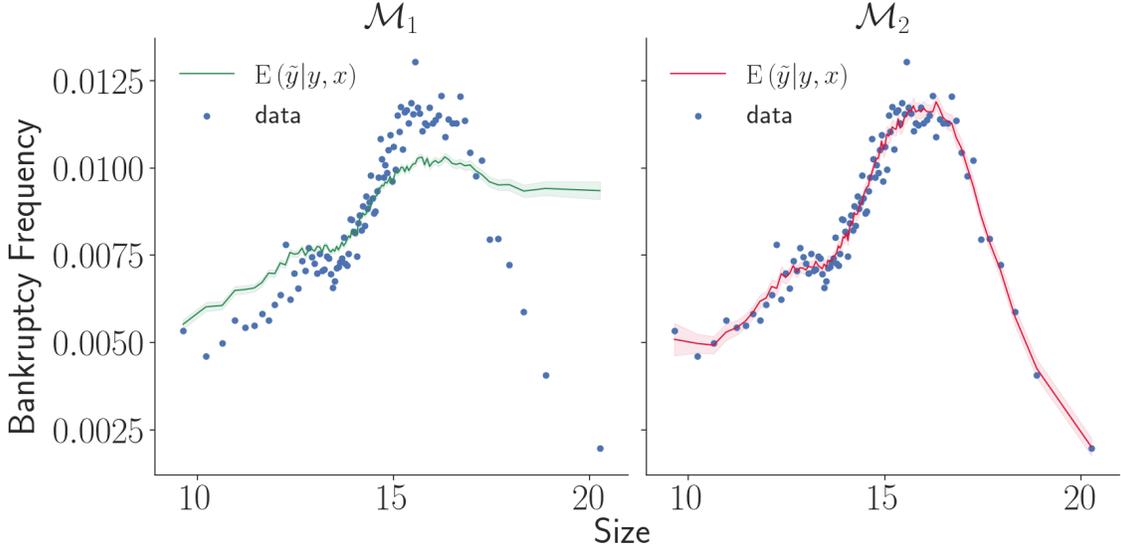


Figure 9: Realized and estimated bankruptcy probabilities. The figure shows the results with respect to the size variable (logarithm of deflated sales) for  $\mathcal{M}_1$  (left panel) and  $\mathcal{M}_2$  (right panel). The data are divided into 100 equally sized groups based on the size variable. For each group, the empirical estimate of the bankruptcy probability is the fraction of bankrupt firms. These empirical estimates are represented as dots, where the corresponding  $x$ -value (size) has been set to the mean within the group. The model estimates for each of the 100 groups are obtained by, for each posterior sample  $\theta$ , averaging the posterior predictive  $\Pr(\tilde{y}_k = 1|\mathbf{y}, x_k)$  for all observations  $k$  in a group, and subsequently computing the posterior predictive mean  $\mathbb{E}(\tilde{y}_k = 1|\mathbf{y}, x_k)$  (solid line) and 90% prediction interval (quantiles 5-95, shaded region).

## 6 Conclusions

A simple and effective approach is proposed to speed up sequential Monte Carlo for static Bayesian models using data subsampling. Its key ingredients are an efficient annealed likelihood estimator and an effective Markov kernel move step based on Hamiltonian Monte Carlo to boost particle diversity. This kernel is computationally expensive for large datasets and data subsampling is crucial to obtain a feasible approach. We argue that the subsampling approach is also very convenient for managing computer memory when implementing SMC using parallel computing, because it avoids the need for each worker to store the full dataset. We demonstrate that the method performs efficiently and accurately for four generalized linear models and a generalized additive model. Moreover, it allows Bayesian model selection through

accurate estimates of the marginal likelihood, which is a major advantage compared to Subsampling MCMC. We also illustrate that the limitation of our method is that its performance depends on good control variates, which can be challenging to construct in certain models. An anonymous reviewer suggested we may use the SMC particles to construct a surrogate function to use as control variate in more complex models. How to do this in a computationally efficient way is an open question, and we leave this extension for future research

## Acknowledgements

We thank the Associate Editor and two reviewers for helping to improve both the content and the presentation of the article. Khue-Dung Dang, David Gunawan, Matias Quiroz and Robert Kohn were partially supported by Australian Research Council Center of Excellence grant CE140100049.

## References

- Baldi, P., Sadowski, P., and Whiteson, D. (2014). Searching for exotic particle in high energy physics with deep learning. *Nature Communications*, 5.
- Bardenet, R., Doucet, A., and Holmes, C. (2017). On Markov chain Monte Carlo methods for tall data. *The Journal of Machine Learning Research*, 18(1):1515–1557.
- Beskos, A., Jasra, A., Kantas, N., and Thiery, A. (2016). On the convergence of adaptive sequential Monte Carlo methods. *The Annals of Applied Probability*, 26(2):1111–1146.
- Betancourt, M. (2017). A conceptual introduction to Hamiltonian Monte Carlo. *arXiv preprint arXiv:1701.02434*.
- Brooks, S., Gelman, A., Jones, G., and Meng, X.-L. (2011). *Handbook of Markov chain Monte Carlo*. CRC press.
- Buchholz, A., Chopin, N., and Jacob, P. E. (2018). Adaptive tuning of Hamiltonian Monte Carlo within sequential Monte Carlo. *arXiv preprint arXiv:1808.07730*.
- Ceperley, D. and Dewing, M. (1999). The penalty method for random walks with uncertain energies. *The Journal of Chemical Physics*, 110(20):9812–9820.

- Chib, S. and Jeliazkov, I. (2001). Marginal likelihood from the Metropolis-Hastings output. *Journal of American Statistical Association*, 96(453):270–281.
- Chopin, N. (2002). A sequential particle filter method for static models. *Biometrika*, 89(3):539–552.
- Dang, K.-D., Quiroz, M., Kohn, R., Tran, M.-N., and Villani, M. (2019). Hamiltonian monte carlo with energy conserving subsampling. *Journal of Machine Learning Research*, 20(100):1–31.
- Daviet, R. (2016). Inference with Hamiltonian sequential Monte Carlo simulators. <http://www.remidaviet.com/files/HSMC-paper.pdf>.
- Del Moral, P., Doucet, A., and Jasra, A. (2006). Sequential Monte Carlo samplers. *Journal of the Royal Statistical Society, Series B*, 68(3):411–436.
- Del Moral, P., Doucet, A., and Jasra, A. (2012). An adaptive Sequential Monte Carlo for approximate Bayesian computation. *Statistics and Computing*, 22(5):1009–1020.
- Deligiannidis, G., Doucet, A., and Pitt, M. K. (2018). The correlated pseudomarginal method. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 80(5):839–870.
- Doucet, A., De Freitas, N., and Gordon, N. (2001). An introduction to sequential Monte Carlo methods. In *Sequential Monte Carlo methods in practice*, pages 3–14. Springer.
- Duan, J. C. and Fulop, A. (2015). Density-tempered marginalised sequential Monte Carlo samplers. *Journal of Business and Economics Statistics*, 33(2):192–202.
- Duane, S., Kennedy, A. D., Pendleton, B. J., and Roweth, D. (1987). Hybrid Monte Carlo. *Physics Letters B*, 195(2):216–222.
- Fearnhead, P. and Taylor, B. M. (2013). An adaptive sequential Monte Carlo sampler. *Bayesian Analysis*, 8(2):411–438.
- Giordani, P., Jacobson, T., Von Schedvin, E., and Villani, M. (2014). Taking the twists into account: Predicting firm bankruptcy risk with splines of financial ratios. *Journal of Financial and Quantitative Analysis*, 49(4):1071–1099.
- Guldás, H., Cemgil, A. T., Whiteley, N., and Heine, K. (2015). A practical introduction to butterfly and adaptive resampling in sequential monte carlo. *IFAC-PapersOnLine*, 48(28):787–792.

- Heine, K., Whiteley, N., and Cemgil, A. T. (2019). Parallelizing particle filters with butterfly interactions. *Scandinavian Journal of Statistics*.
- Jasra, A., Stephens, D. A., Doucet, A., and Tsagaris, T. (2011). Inference for Lévy-driven stochastic volatility models via adaptive Sequential Monte Carlo. *Scandinavian Journal of Statistics*, 38(1):1–22.
- Jeffreys, H. (1961). *The Theory of Probability*. OUP Oxford.
- Johnson, A. A., Jones, G. L., and Neath, R. C. (2013). Component-wise Markov chain Monte Carlo: Uniform and geometric ergodicity under mixing and composition. *Statistical Science*, 28(3):360–375.
- Kass, R. E. and Raftery, A. E. (1995). Bayes factors. *Journal of American Statistical Association*, 90(430):773–795.
- Lee, A., Yau, C., Giles, M. B., Doucet, A., and Holmes, C. C. (2010). On the utility of graphics cards to perform massively parallel simulation of advanced monte carlo methods. *Journal of computational and graphical statistics*, 19(4):769–789.
- Liu, J. S. (2001). *Monte Carlo strategies in scientific computing*. New York: Springer.
- Murray, L. M., Lee, A., and Jacob, P. E. (2016). Parallel resampling in the particle filter. *Journal of Computational and Graphical Statistics*, 25(3):789–805.
- Neal, R. (2001). Annealed importance sampling. *Statistics and Computing*, 11:125–139.
- Neal, R. M. (2011). MCMC using Hamiltonian dynamics. *Handbook of Markov chain Monte Carlo*.
- Quiroz, M., Kohn, R., Villani, M., and Tran, M. N. (2019). Speeding up MCMC by efficient data subsampling. *Journal of American Statistical Association*, 114:831–843.
- Quiroz, M., Tran, M.-N., Villani, M., Kohn, R., and Dang, K.-D. (2018a). The block-Poisson estimator for optimally tuned exact subsampling MCMC. *arXiv preprint arXiv:1603.08232v5*.
- Quiroz, M. and Villani, M. (2013). Dynamic mixture-of-experts models for longitudinal and discrete-time survival data. <https://github.com/mattiasvillani/Papers/raw/master/DynamicMixture.pdf>.

- Quiroz, M., Villani, M., Kohn, R., Tran, M.-N., and Dang, K.-D. (2018b). Subsampling MCMC: An introduction for the survey statistician. *Sankhya A*, 80.
- Roberts, G. O., Gelman, A., and Gilks, W. R. (1997). Weak convergence and optimal scaling of random walk Metropolis-Hastings. *Annals of Applied Probability*, 7(1):110–120.
- Roberts, G. O. and Stramer, O. (2002). Langevin diffusions and Metropolis-Hastings algorithms. *Methodology and Computing in Applied Probability*, 4(4):337–357.
- Sim, A., Filippi, S., and Stumpf, M. P. (2012). Information geometry and sequential Monte Carlo. *arXiv preprint arXiv:1212.0764*.
- South, L. F., Pettitt, A. N., and Drovandi, C. C. (2016). Sequential Monte Carlo for static Bayesian models with independent MCMC proposals. <https://core.ac.uk/download/pdf/78105120.pdf>.
- South, L. F., Pettitt, A. N., Friel, N., and Drovandi, C. C. (2017). Efficient use of derivative information within SMC methods for static Bayesian Models. <https://eprints.qut.edu.au/108150/>.
- Tran, M. N., Kohn, R., Quiroz, M., and Villani, M. (2017). The block-pseudo marginal sampler. *preprint arXiv:1603.02485v5*.
- Wang, L., Wang, S., and Bouchard-Côté, A. (2019). An annealed sequential Monte Carlo method for Bayesian phylogenetics. *arXiv preprint arXiv:1806.08813v3*.