

Série des Documents de Travail

n° 2017-40

**Informed Sub-Sampling MCMC: Approximate
Bayesian Inference for Large Datasets**

F. MAIRE¹

N. FRIEL²

P. ALQUIER³

Les documents de travail ne reflètent pas la position du CREST et n'engagent que leurs auteurs.
Working papers do not reflect the position of CREST but only the views of the authors.

¹ School of Mathematics and Statistics, University College Dublin; Insight Centre for Data Analytics, University College Dublin. E-mail : florian.maire@ucd.ie

² School of Mathematics and Statistics, University College Dublin; Insight Centre for Data Analytics, University College Dublin. E-mail : nial.friel@ucd.ie

³ CREST-ENSAE, France. E-mail : Pierre.alquier@ensae.fr

Informed Sub-Sampling MCMC: Approximate Bayesian Inference for Large Datasets

Florian Maire^{a,b,*}, Nial Friel^{a,b}, Pierre Alquier^c

^a*School of Mathematics and Statistics, University College Dublin*

^b*Insight Centre for Data Analytics, University College Dublin*

^c*CREST, ENSAE, Université Paris Saclay*

Abstract

This paper introduces a framework for speeding up Bayesian inference conducted in presence of large datasets. We design a Markov chain whose transition kernel uses an unknown fraction of fixed size of the available data that is randomly refreshed throughout the algorithm. Inspired by the Approximate Bayesian Computation (ABC) literature, the subsampling process is guided by the fidelity to the observed data, as measured by summary statistics. The resulting algorithm, Informed Sub-Sampling MCMC, is a generic and flexible approach which, contrarily to existing scalable methodologies, preserves the simplicity of the Metropolis–Hastings algorithm. Even though exactness is lost, *i.e* the chain distribution approximates the target, we study and quantify theoretically this bias and show on a diverse set of examples that it yields excellent performances when the computational budget is limited. If available and cheap to compute, we show that setting the summary statistics as the maximum likelihood estimator is supported by theoretical arguments.

Keywords: Bayesian inference, Big-data, Approximate Bayesian Computation, noisy Markov chain Monte Carlo

Primary: 65C40, 65C60 – Secondary: 62F15

1. Introduction

The development of statistical methodologies that scale to large datasets represents a significant research frontier in modern statistics. This paper presents a generic and flexible approach to directly address this challenge when a Bayesian strategy is followed. Given a set of observed data (Y_1, \dots, Y_N) , a specified prior distribution p and a likelihood function f , estimating parameters $\theta \in \Theta$ of the model proceeds via exploration of the posterior distribution π defined on $(\Theta, \mathcal{B}(\Theta))$ by

$$\pi(d\theta | Y_1, \dots, Y_N) \propto f(Y_1, \dots, Y_N | \theta)p(d\theta). \quad (1)$$

*Corresponding author

Email address: florian.maire@ucd.ie (Florian Maire)

Stochastic computation methods such as Monte Carlo methods allow one to estimate characteristics of π . In Bayesian inference, Markov chain Monte Carlo (MCMC) methods remain the most widely used strategy. Paradoxically, improvements in data acquisition technologies together with increased storage capacities, present a new challenge for these methods. Indeed, the size of the data set N (along with the dimension of each observation) can become so large, that even a routine likelihood evaluation is made prohibitively computationally intensive. As a consequence, MCMC methods such as the Metropolis–Hastings algorithm (Metropolis et al., 1953) cannot be considered for reasonable runtime. This issue has recently generated a lot of research activity, see Bardenet et al. (2015) for a comprehensive review.

Most of the scalable MCMC methods proposed in the literature are based on approximations of the Metropolis-Hastings (M-H) algorithm. In the sequel, we will refer to as *exact approximations*, algorithms that produce samples from the target distribution when the chain is in stationary regime, as opposed to *approximate* methods that do not. Central to those scalable MCMC approaches is the idea that only the calculation of the likelihood of a subset of data would be required to simulate a new state of the Markov chain. Following the development of pseudo-marginal algorithms (Andrieu and Roberts, 2009; Andrieu and Vihola, 2015), a first direction has been to replace the likelihoods in the M-H acceptance ratio by positive unbiased estimators (based on a subset of data). Although appealing since exact, this approach remains (for now) mostly theoretical because such estimators are in general not available (Jacob et al., 2015). Attempts to circumvent the positivity and unbiasedness requirements of the estimator have been studied in Quiroz et al. (2016) and Quiroz et al. (2015) respectively. In both cases, the authors resort to sophisticated control variates, which can be computationally expensive to compute.

Other authors have proposed to approximate the log-likelihood ratio by subsampling data points (Korattikara et al. (2014); Bardenet et al. (2014, 2015)), the objective being to mimic the accept/reject decision that would be achieved by the Metropolis-Hastings algorithm. Even though the resulting algorithms are not exact, the *Confidence sampler* proposed in Bardenet et al. (2014) and refined in Bardenet et al. (2015) is designed such that the accept/reject decision is, with an arbitrarily high probability, identical to that taken by the Metropolis-Hastings algorithm. The construction of this algorithm, based on concentration inequalities, allows to bound the L1 distance between the stationary distribution of the algorithm and π . The price to pay is that the number of likelihood evaluations is not fixed but adaptively set by the algorithm at each iteration and, as noted in Bardenet et al. (2015), it is of order $\mathcal{O}(N)$ when the chain reaches equilibrium. This number can be brought down if an accurate proxy of the log-likelihood ratio, acting as control variate, is available, as demonstrated in Bardenet et al. (2015).

More recently, a stream of research has shed light on the use of continuous time Markov processes (Zig-Zag process, Langevin diffusion) to perform Bayesian analysis of tall dataset (Bierkens et al., 2016; Pollock et al., 2016; Fearnhead et al., 2016). The computational bottleneck for this class of methods is the calculation of the gradient of the log-likelihood and it has been shown that provided that an unbiased estimate is used, they remain exact. Here again, the use of control variates to reduce the variance of the estimator is in practice essential to reach the full potential of these methods. However, we note that those approaches represent a significant departure from the M-H algorithm and as such lose its implementational simplicity.

In this paper, we propose Informed Sub-Sampling MCMC, a novel methodology which aims to make the best use of a computational resource available for a given computational run-time, while still preserving the celebrated simplicity of the standard M–H sampler. The state space Θ is extended with a n -dimensional vector of unique integers $U_k \subset \{1, \dots, N\}$ identifying a subset of the data used by the Markov transition kernel at the k -th iteration of the algorithm, where $n \ll N$ is set according to the available computational budget. Central to our approach is the fact that each subset is weighted according to a *similarity measure* with respect to the full set of data through summary statistics, in the spirit of Approximate Bayesian Computation (ABC) (see *e.g.* [Marin et al., 2012](#)). The subset variable is randomly refreshed at each iteration according to the similarity measure. The Markov chain transition kernel only uses a fraction n/N of the available data which is by construction –and contrary to [Maclaurin and Adams \(2015\)](#), [Korattikara et al. \(2014\)](#) and [Bardenet et al. \(2014\)](#)– held constant throughout the algorithm. Moreover, unlike most of the papers mentioned before, our method can be applied to virtually any model (involving *i.i.d.* data or not), as it does not require any assumption on the likelihood function nor on the prior distribution. Our algorithm can be cast as a *noisy* MCMC method since the marginal in θ of our Markov chain targets an approximation of π that we quantify using the framework established in [Alquier et al. \(2016\)](#). In the special case where the data are *i.i.d.* realizations from an exponential model, we prove that when the summary statistics is set as the sufficient statistics, this yields an optimal approximation, in the sense of minimizing an upper bound of the Kullback-Leibler (KL) divergence between π and the marginal target of our method. In the general case, we show that setting the summary statistics as the maximum likelihood estimator allows to bound the approximation error (in L1 distance) of our algorithm. We connect our work to a number of recent papers including [Rudolf and Schweizer \(2015\)](#); [Huggins and Zou \(2016\)](#); [Dalalyan \(2017\)](#) that bound approximation error of MCMC algorithms, using the Wasserstein metric.

To summarize, the main contribution of our work is to show that, under verifiable conditions, it is possible to infer π through a scalable approximation of the M–H algorithm where the computational budget of each iteration is fixed (through the subset size n). To do so, it is necessary to draw the subsets according to a similarity measure with respect to the full data set and not uniformly at random, as previously explored in the literature. We show that setting the similarity measure as the squared L2 distance between the full dataset and subsample maximum likelihood estimators is supported by theoretical arguments.

Section 2 presents a striking real data example which we hope will help the reader to understand the problem we address and motivate the solution we propose, without going into further technical details at this stage. In Section 3, we provide theoretical results concerning exponential-family models, which we illustrate through a probit example. This section allows us to justify our motivations supporting the Informed Sub-Sampling general methodology which is rigorously presented in Section 4. In Section 5, we study the transition kernel of our algorithm and show that it yields a Markov chain targeting, marginally, an approximation of π . The approximation error is quantified and we provide theoretical justifications for setting up the Informed Sub-Sampling tuning parameters, including the choice of summary statistics. Finally, in Section 6, our method is used to estimate parameters of an autoregressive time series and a logistic regression model. It is also illustrated to perform a binary classification task. In the latter example, we

compare the performance of our algorithm with the SubLikelihoods approach proposed in [Bardinet et al. \(2014\)](#).

2. An introductory example

We showcase the principles of our approach on a first real data example. The problem at hand is to infer some template shapes of handwritten digits from the MNIST database (<http://yann.lecun.com/exdb/mnist/>).

Example 1. *The data Y_1, Y_2, \dots are modelled by a deformable template model ([Allasonnière et al., 2007](#)). Each data Y_i is a 15×15 pixel image representing an handwritten digit whose conditional distribution given its class $J(i) \in (0, 1, \dots, 9)$ is a random deformation of the template shape, parameterized by a $d = 256$ dimensional vector $\theta_{J(i)}$. Assuming small deformations, the model is similar to a standard regression problem:*

$$Y_i = \phi(\theta_{J(i)}) + \sigma^2 \epsilon_i, \quad (2)$$

where Y_i is regarded as a vector \mathbb{R}^{225} , $\phi : \mathbb{R}^{256} \rightarrow \mathbb{R}^{225}$ is some deterministic mapping and $\sigma > 0$ is the standard deviation of the additive noise $\epsilon_i \sim \mathcal{N}(0_{225}, \text{Id}_{225})$.

Given a set of N labeled images Y_1, Y_2, \dots, Y_N and a prior distribution for $\theta = \{\theta_1, \dots, \theta_9\}$, one can estimate θ through its posterior distribution π , for example using the Metropolis–Hastings (M-H) algorithm ([Metropolis et al., 1953](#)). However, since the regression function ϕ in (2) is quite sophisticated, even a single likelihood evaluation is expensive to calculate. As a result, the M-H efficiency can be questioned as computing the N likelihoods in the M-H ratio dramatically slows down each transition.

At this stage, we do not provide precise details on the Informed Sub-Sampling MCMC method but we simply provide an insight of the rationale of our approach. It designs a Markov chain whose transition kernel targets a scaled version of the posterior distribution of the parameter of interest θ given a random subset of n images ($n \ll N$). More specifically, we inject in the standard M–H transition a decision about *refreshing* the subset of data, which, as a result, will change randomly over time. In this example, we use the knowledge of the observation labels to promote subsets of images in which the proportion of each digit is balanced.

We consider $N = 10,000$ images of five digits $1, \dots, 5$, subsets of size $n = 100$ and a non-informative Gaussian prior for θ , as specified in [Allasonnière et al. \(2007\)](#). Figure 1 indicates a striking advantage of our method compared to a standard M-H using the same $N = 10,000$ images. In this scenario, we allow a fixed computational budget (1 hour) for both methods and compare the estimation of the mean estimate of the two Markov chains. Qualitatively, the upper part of Figure 1 compares the estimated template shapes of the five digits at different time steps and shows that our method allows to extract template shapes much quicker than the standard M-H, while still reaching an apparent similar graphical quality after one hour. This fact is confirmed quantitatively, in the lower part of Figure 1, which plots, against time and for both methods, the Euclidean distance between the Markov chain mean estimate and the maximum likelihood estimate ($\theta_1^*, \dots, \theta_5^*$) obtained using a stochastic EM ([Allasonnière et al., 2007](#)). More precisely,

we compare the real valued function $\{d(t), t \in \mathbb{R}\}$ defined as

$$d(t) = \sum_{j=1}^5 \|\theta_j^* - \mu(\theta_{j,1:\kappa(t)})\|, \quad \text{where} \quad \begin{cases} \forall t \in \mathbb{R}, \kappa(t) = \max_{k \in \mathbb{N}} \{t \geq \tau_k\}, \\ \tau_k \text{ is the time at the end of the } k\text{-th iteration,} \\ \forall k \in \mathbb{N}, \mu(\theta_{j,1:k}) = (1/k) \sum_{\ell=1}^k \theta_{j,\ell}, \end{cases}$$

where we have defined for $(j, k) \in \{1, \dots, 5\} \times \mathbb{N}$, $\theta_{j,k}$ as the j -th class parameter obtained after k iterations of the Markov chains.

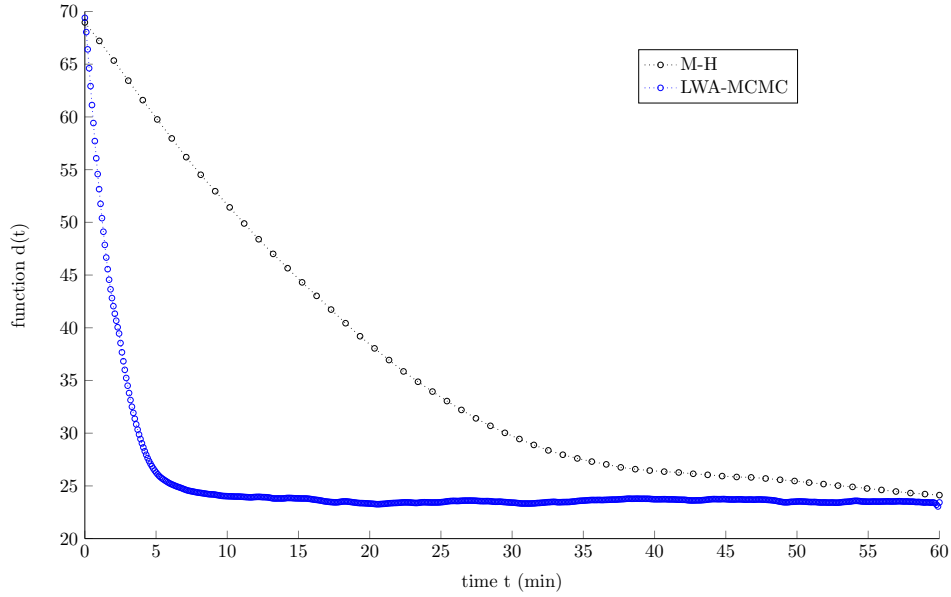
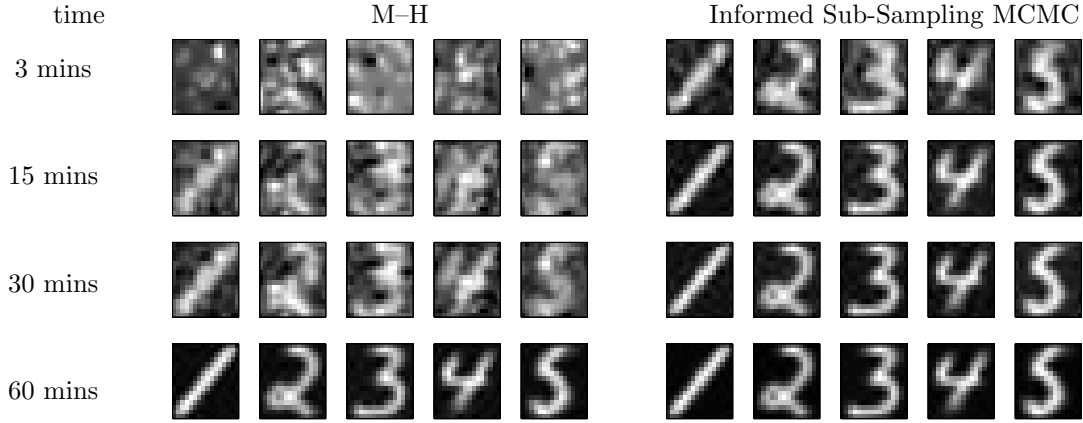


Figure 1: (Example 1: Handwritten digits) Efficiency of template estimation through M-H and Informed Sub-Sampling MCMC.

One can see that the transient phase of the Informed Sub-Sampling Markov chain is significantly shorter than that of the Metropolis-Hastings chain. More details on

Example 1 can be found at Section 6. In particular, Figure 13 shows that the stationary distribution of Informed Sub-Sampling matches reasonably well π , which is a primary concern for Bayesian inference.

Our algorithm provides very encouraging results for this real data example. We motivate and formalize our method in Sections 3 and 4 and provide theoretical arguments supporting it at Section 5.

3. Approximation of the posterior distribution in exponential models: an optimality result

In this section, we consider the case of N independent and identically distributed (*i.i.d.*) observations from an exponential model. Sampling from the posterior distribution of such models using the Metropolis-Hastings algorithm is effortless since the information conveyed by the N observations is contained in the sufficient statistics vector, which needs to be calculated only once.

The existence of sufficient statistics in this type of models allows us to establish a number of theoretical results that will be used to design and justify our Informed Sub-Sampling methodology that approximately samples from posterior distributions in general contexts, *i.e.* non-*i.i.d.* observations from general likelihood models without sufficient statistics. More precisely, Propositions 1 and 2 put forward an optimal approximation of the posterior distribution π by a distribution $\tilde{\pi}_n$ of the parameter of interest given only a subsample of n observations. Finally, Proposition 3 justifies the introduction of a probability distribution on the set of subsamples. This is an essential element of our work as it represents a significant departure compared to all existing subsampling methodologies proposed in the Markov chain Monte Carlo literature, that have assumed uniform distribution on the subsamples.

3.1. Notation

Let $(Y_1, \dots, Y_N) \in \mathcal{Y}^N$ be a set of *i.i.d.* observed data ($\mathcal{Y} \subseteq \mathbb{R}^m$, $m > 0$) and define

- $Y_{i:j} = (Y_i, \dots, Y_j)$ if $1 \leq i \leq j \leq N$ with the convention that $Y_{i:i} = \{Y_i\}$, otherwise.
- $Y_U = \{Y_k, k \in U\}$, where $U \subseteq \{1, \dots, N\}$.

In this section, we assume that the likelihood model f belongs to the exponential family and is fully specified by a vector of parameters $\theta \in \Theta$, ($\Theta \subseteq \mathbb{R}^d$, $d > 0$), a bounded mapping $g : \Theta \rightarrow \mathcal{S}$ and a sufficient statistic mapping $S : \mathcal{Y} \rightarrow \mathcal{S}$ ($\mathcal{S} \subseteq \mathbb{R}^s$, $s > 0$) such that

$$f(y|\theta) = \exp\{g(\theta)^T S(y)\} / L(\theta), \quad L(\theta) = \int_{\mathcal{Y} \in \mathcal{Y}} \exp\{S(y)^T g(\theta)\} dy,$$

is the density of the likelihood distribution with respect to the Lebesgue measure. The posterior distribution π is defined on the measurable space $(\Theta, \mathcal{B}(\Theta))$ by its density function

$$\pi(\theta | Y_{1:N}) = p(\theta) \frac{\exp\left\{\sum_{k=1}^N S(Y_k)^T g(\theta)\right\}}{L(\theta)^N} / Z(Y_{1:N}), \quad (3)$$

where

$$Z(Y_{1:N}) = \int p(d\theta) \frac{\exp\left\{\sum_{k=1}^N S(Y_k)^T g(\theta)\right\}}{L(\theta)^N}. \quad (4)$$

p is a prior distribution defined on $(\Theta, \mathcal{B}(\Theta))$ and with some abuse of notation, p denotes also the probability density function on Θ .

For all $n \leq N$, we define \mathbf{U}_n as the set of the possible combinations of n different integer numbers less than or equal to N and \mathcal{U}_n as the powerset of \mathbf{U}_n . In the sequel, we set n as a constant and wish to compare the posterior distribution π (3) with any distribution from the family $\mathbf{F}_n = \{\tilde{\pi}_n(U), U \in \mathbf{U}_n\}$, where for all $U \in \mathbf{U}_n$, we have defined $\tilde{\pi}_n(U)$ as the distribution on $(\Theta, \mathcal{B}(\Theta))$ with probability density function

$$\tilde{\pi}_n(\theta | Y_U) \propto p(\theta) f(Y_U | \theta)^{N/n}. \quad (5)$$

3.2. Optimal subsets for the Kullback-Leibler divergence between π and $\tilde{\pi}_n$

Recall that for two measures π and $\tilde{\pi}$ defined on the same measurable space $(\Theta, \mathcal{B}(\Theta))$, the Kullback-Leibler (KL) divergence between π and $\tilde{\pi}$ is defined as:

$$\text{KL}(\pi, \tilde{\pi}) = \mathbb{E}_\pi \left\{ \log \frac{\pi(\theta)}{\tilde{\pi}(\theta)} \right\}. \quad (6)$$

Although not a proper distance between probability measures defined on the same space, $\text{KL}(\pi, \tilde{\pi})$ is used as a similarity criterion between π and $\tilde{\pi}$. It can be interpreted in information theory as a measure of the information lost when $\tilde{\pi}$ is used to approximate π , which is our primary concern here. We now state the main result of this section.

Proposition 1. *For any subset $U \in \mathbf{U}_n$, define the vector of difference of sufficient statistics between the whole dataset and the subset Y_U as*

$$\Delta_n(U) = \sum_{k=1}^N S(Y_k) - (N/n) \sum_{k \in U} S(Y_k). \quad (7)$$

Then, the following inequality holds:

$$\text{KL} \{ \pi, \tilde{\pi}_n(U) \} \leq B(Y, U), \quad (8)$$

where

$$B(Y, U) = \log \mathbb{E}_\pi \exp \{ \| \mathbb{E}_\pi(g(\theta)) - g(\theta) \| \| \Delta_n(U) \| \}. \quad (9)$$

The proof is postponed to [Appendix A.1](#) and follows from straightforward algebra and applying Cauchy-Schwartz inequality. The following corollary is an immediate consequence of Proposition 1.

Corollary 1. (i) *Define the set:*

$$\mathbf{U}_n^* := \left\{ U \in \mathbf{U}_n, \quad \frac{1}{N} \sum_{k=1}^N S(Y_k) = \frac{1}{n} \sum_{k \in U} S(Y_k) \right\}. \quad (10)$$

If \mathbf{U}_n^ is non-empty, then for any $U \in \mathbf{U}_n^*$, then $\pi(\theta | Y) = \tilde{\pi}_n(\theta | Y_U)$, π -almost everywhere.*

(ii) Let $(U_1, U_2) \in \mathcal{U}_n^2$. Assume $\|\Delta_n(U_1)\| \leq \|\Delta_n(U_2)\|$, then $B(Y, U_1) \leq B(Y, U_2)$.

A stronger result can be obtained under the assumption that a Bernstein-von Mises Theorem [Van der Vaart \(2000\)](#); [Le Cam \(1986\)](#) holds for the concentration of π to its Normal approximation:

$$\hat{\pi}(\cdot | Y_{1:N}) := \mathcal{N}(\theta^*(Y_{1:N}), I^{-1}(\theta_0)/N), \quad (11)$$

where \mathcal{N} denotes the Normal distribution, $\theta^*(Y_{1:N}) = \arg \max_{\theta \in \Theta} f(Y_{1:N} | \theta)$, $\theta_0 \in \Theta$ is some parameter and $I(\theta)$ is the Fisher information matrix given $Y_{1:N}$ at θ .

Proposition 2. Let $(U_1, U_2) \in \mathcal{U}_n^2$. Assume that for all $i \in \{1, \dots, d\}$, $|\Delta_n(U_1)^{(i)}| \leq |\Delta_n(U_2)^{(i)}|$, where $|\Delta_n(U_1)^{(i)}|$ refers to the i -th element of $\Delta_n(U_1)$ (7). Then $\widehat{\text{KL}}_n(U_1) \leq \widehat{\text{KL}}_n(U_2)$, where $\widehat{\text{KL}}_n(U)$ is the Kullback-Leibler divergence between the asymptotic approximation of the posterior $\hat{\pi}$ (11) and $\tilde{\pi}_n(U)$ (5).

The proof is postponed to [Appendix A.2](#). Note that the asymptotic approximation is for $N \rightarrow \infty$ and for a fixed n and is thus relevant to the context of our analysis.

3.3. Weighting the subsamples

Let us consider the distribution $\nu_{n,\epsilon}$ on the discrete space \mathcal{U}_n defined for all $\epsilon \geq 0$ by:

$$\nu_{n,\epsilon}(U) \propto \exp\{-\epsilon \|\Delta_n(U)\|^2\}, \quad \text{for all } U \in \mathcal{U}_n. \quad (12)$$

$\nu_{n,\epsilon}$ assigns a weight to any subset according to their representativeness with respect to the full dataset. When $\epsilon = 0$, $\nu_{n,\epsilon}$ is uniform on \mathcal{U}_n while when $\epsilon \rightarrow \infty$, $\nu_{n,\epsilon}$ is uniform on the set of subset(s) that minimize(s) $U \mapsto \|\Delta_n(U)\|$. [Proposition 2](#) suggests that for exponential models, the optimal inference based on subsamples of size n is obtained by picking the subposterior $\pi_n(U)$ (5) using the distribution $U \sim \nu_{n,\epsilon}$ with $\epsilon \rightarrow \infty$.

We now state [Proposition 3](#). This result is important even though somewhat obscure at this stage. Indeed, we will show that it is a necessary condition for the method we introduce in [Section 4](#) to converge. In fact, moving away to general models (*i.e.* non *i.i.d.* and non exponential) amounts to relax the sufficient statistics existence assumption as well as the $\epsilon \rightarrow \infty$ condition. This will be achieved by constructing a class of summary statistics for the model at hand for which a similar result to [Proposition 3](#) holds.

Proposition 3. For any $\theta \in \Theta$ and $\epsilon > 0$, there exists $M < \infty$ such that:

$$\mathbb{E}_\epsilon \left\{ \frac{f(Y | \theta)}{f(Y_U | \theta)^{N/n}} \right\} < M, \quad (13)$$

where \mathbb{E}_ϵ is the expectation under $\nu_{n,\epsilon}$, as defined in (12).

The proof is postponed to [Appendix A.3](#). Note that [Proposition 3](#) essentially holds because $\log \nu_{n,\epsilon}$ is quadratic in $\|\Delta_n(U)\|$. Other weighting schemes for the subsets (e.g. uniform weights or weights $\propto \exp\{\epsilon \|\Delta_n(U)\|\}$) would not necessarily allow to bound $\mathbb{E}\{f(Y | \theta)/f(Y_U | \theta)^{N/n}\}$.

n	$\ \Delta_n(U)\ $	KL $\{\pi, \tilde{\pi}_n(U)\}$	$B(Y, U)$
1,000	3	0.004	0.04
1,000	14	0.11	0.18
1,000	23	0.19	0.29
100	33	0.41	0.54

Table 1: (Example 2: Probit model) Comparison of the KL divergence between π and the optimal $\tilde{\pi}_n \in \mathbf{F}_n$ ($\|\Delta_n(U)\| = 3$) and other distributions in \mathbf{F}_n .

3.4. Illustration with a probit model: effect of choice of sub-sample

We consider a pedagogical example, based on a probit model, to illustrate the results from the previous subsections.

Example 2. A probit model is used in regression problems in which a binary variable $Y_k \in \{0, 1\}$ is observed through the following sequence of independent random experiments, defined for all $k \in \{1, \dots, N\}$ as:

(i) Draw $X_k \sim \mathcal{N}(\theta^*, \gamma^2)$

(ii) Set Y_k as follows

$$Y_k = \begin{cases} 1, & \text{if } X_k > 0, \\ 0, & \text{otherwise.} \end{cases} \quad (14)$$

Observing a large number of realizations Y_1, \dots, Y_N , we aim to estimate the posterior distribution of θ . If γ is unknown, the model is not identifiable and for simplicity we considered it as known here. The likelihood function can be expressed as

$$f(Y_k | \theta) = \alpha(\theta)^{Y_k} (1 - \alpha(\theta))^{(1-Y_k)} = (1 - \alpha(\theta)) \left(\frac{\alpha(\theta)}{1 - \alpha(\theta)} \right)^{Y_k}, \quad (15)$$

where $\alpha(\theta) = \int_0^\infty (2\pi\gamma^2)^{-1/2} \exp\{-(1/2\gamma^2)(t - \theta)^2\} dt$ and clearly belongs to the exponential family. The pdf of the posterior distribution π and any distribution $\tilde{\pi}_n(U) \in \mathbf{F}_n$ writes respectively as

$$\pi(\theta | Y_{1:N}) \propto p(\theta) (1 - \alpha(\theta))^N \left(\frac{\alpha(\theta)}{1 - \alpha(\theta)} \right)^{\sum_{k=1}^N Y_k},$$

$$\tilde{\pi}_n(\theta | Y_U) \propto p(\theta) (1 - \alpha(\theta))^N \left(\frac{\alpha(\theta)}{1 - \alpha(\theta)} \right)^{(N/n) \sum_{k \in U} Y_k},$$

where p is a prior density on θ . Again, in this example, the posterior density is easy to evaluate pointwise, even when N is extremely large, as it only requires to sum over all the binary variables Y_1, \dots, Y_N . As a consequence, samples from π can routinely be obtained by a standard M–H algorithm and similarly for any distribution $\tilde{\pi}_n(U) \in \mathbf{F}_n$.

We simulated $N = 10,000$ simulated data Y_1, \dots, Y_N from (14), with true parameter $\theta^* = 1$. We used the prior distribution $p = \mathcal{N}(0, 10)$. In this probit model, S is simply the identity function, implying that $\|\Delta_n(U)\|$ gives the absolute value of the difference between the scaled proportion of 1 and 0's between the full dataset and the subset Y_U .

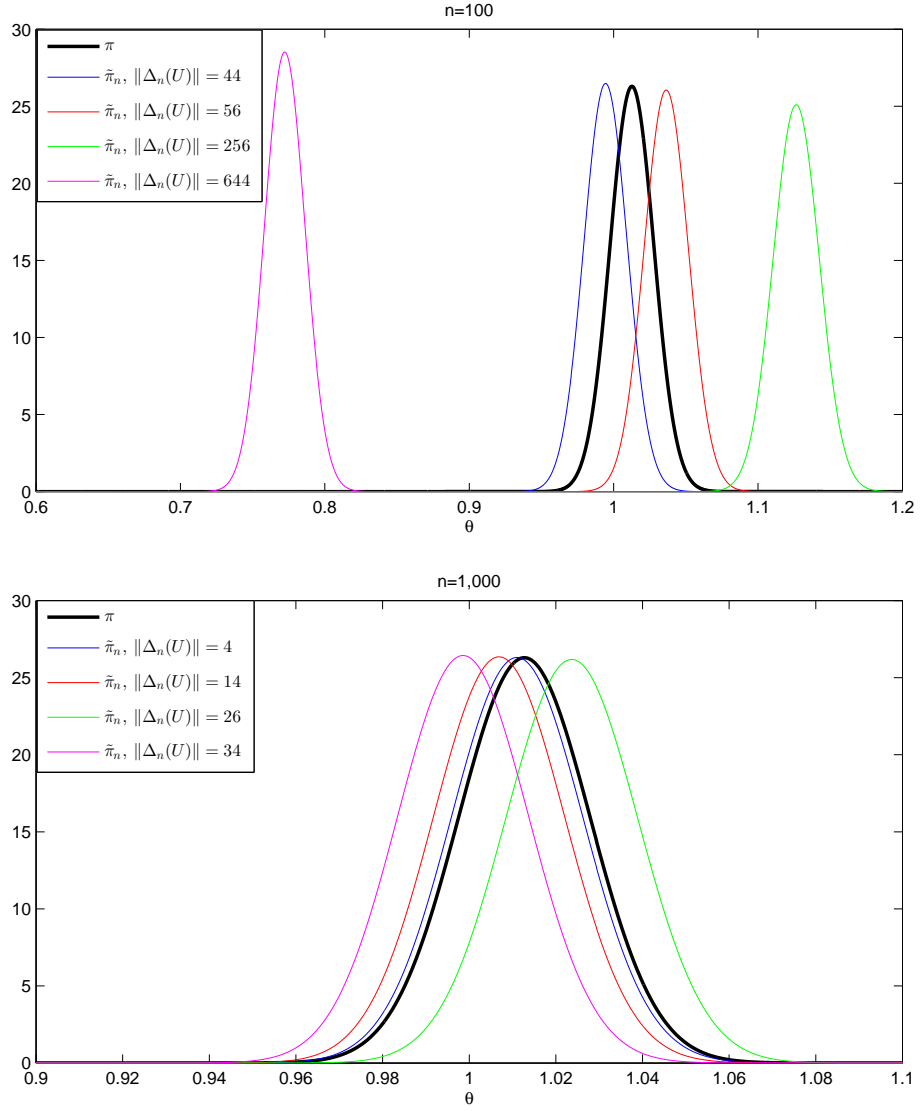


Figure 2: (Example 2: Probit model) Influence of the parameter $U \in \mathcal{U}_n$ on the sub-posterior distribution $\tilde{\pi}_n(U)$ and comparison with π for subsets of size $n = 100$ (top) and $n = 1,000$ (bottom).

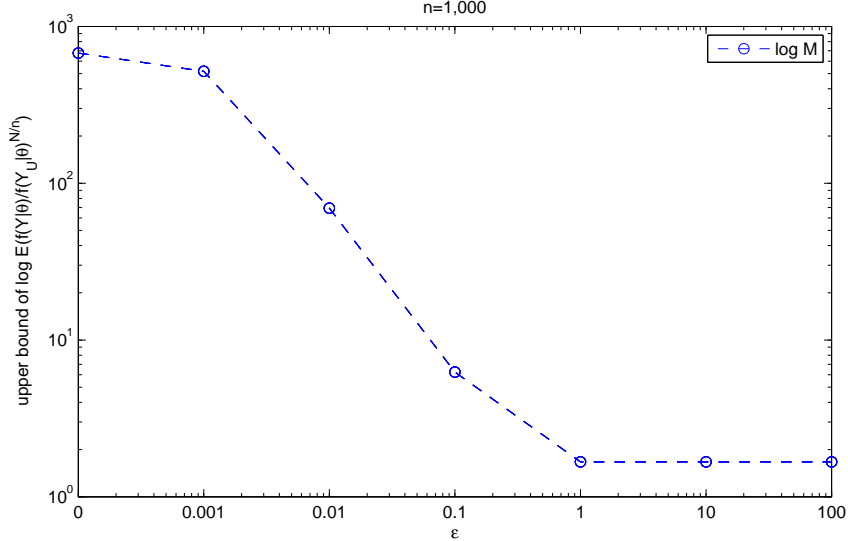


Figure 3: (Example 2: Probit model) Influence of the parameter ϵ of the distribution $\nu_{n,\epsilon}$ on the upper bound M of $\mathbb{E}\{f(Y|\theta)/f(Y_U|\theta)^{N/n}\}$ for $\theta \in (0, 1.5)$, for $n = 1,000$. When $\epsilon = 0$, $\nu_{n,\epsilon}$ is uniform (i.e. an identical weight is assigned to all the subsamples) and as a consequence $M \equiv \infty$. Conversely, when $\epsilon \gg 0$, the mass of $\nu_{n,\epsilon}$ spreads over the best subsamples $Y_U, U \in \mathcal{U}_n$ (i.e. those minimizing $\Delta_n(U)$) and the bound M is smaller than e^2 . Indeed, by assigning a weight $\nu_{n,\epsilon}(U) \propto \exp\{-\epsilon\Delta_n(U)^2\}$ those subsamples $Y_U, U \in \mathcal{U}_n$ that have a large $\Delta_n(U)$ will yield a negligible contribution to the expectation, hence preventing from divergence.

Figure 2 reports the density functions of π and several other distributions $\tilde{\pi}_n(U) \in \mathcal{F}_n$, for $n = 100$ and $n = 1,000$, with different values for the quantity $\|\Delta_n(U)\|$ (7). This plot, as well as the quantitative result of Table 1 are consistent with the statement of Corollary 1: when learning from a subsample of n data, one should work with a subset U featuring a perfect match with the full dataset, i.e. $\|\Delta(U)\| = 0$, or as small as possible to achieve an *optimal* approximation of π . Finally, Figure 3 illustrates Proposition 3: assigning the distribution $\nu_{n,\epsilon}$ (12) to the subsamples allows to control the expectation of the likelihood ratio $f(Y|\theta)/f(Y_U|\theta)^{N/n}$ around 1.

4. Informed Sub-Sampling MCMC

In this section, we do not assume any particular correlation pattern for the sequence of observations, nor any specific likelihood model and simply write the posterior distribution π as

$$\pi(d\theta | Y_{1:N}) \propto p(d\theta) f(Y_{1:N} | \theta). \quad (16)$$

The Informed Sub-Sampling MCMC methodology (IIS-MCMC for short) that we describe now can be regarded as an extension of the approximation detailed in the previous section to non-exponential family models with possibly dependent observations.

4.1. Motivation of our approach

Central to our approach is the idea that all the subsamples Y_U ($U \in \mathcal{U}_n$) are not equally valuable for inferring π . Here, we do not assume the existence of a sufficient statistic mapping for the models under consideration. Thus, in order to discriminate between different subsamples, we introduce an *artificial* summary statistic mapping $S : Y_n \rightarrow \mathcal{S}$ ($n \leq N$), where $\mathcal{S} \subseteq \mathbb{R}^s$. The choice of the summary statistics S is problem specific and is meant to be the counterpart of the sufficient statistic mapping for general models (hence sharing, slightly abusively, the same notation). Since the question of specifying summary statistics also arises in Approximate Bayesian Computation (ABC), one can take advantage of the abundant ABC literature on this topic to find some examples of summary statistics for usual likelihood models (see *e.g.* Nunes and Balding, 2010; Csilléry et al., 2010; Marin et al., 2012; Fearnhead and Prangle, 2012). More details on validation of summary statistics are discussed in Section 5.3.

Because the statistics used to assess the representativeness of a subsample Y_U w.r.t. the full dataset Y are only *summary* and not *sufficient*, the results of Section 3 are no longer valid. In particular, should an optimal subset U^* minimising a distance between $S(Y_U)$ and $S(Y)$ exist, inferring π through the approximation $\tilde{\pi}_n(U^*)$ is in no sense optimal. In fact, as shown in several examples of Section 6, this approximation is usually poor. In such a setting, it is reasonable to consider extending the set of subsamples of interest to a pool of *good* subsamples. This naturally suggests using the distribution $\nu_{n,\epsilon}$ (12) to discriminate between the subsamples, replacing sufficient by summary statistics and relaxing the assumption $\epsilon \rightarrow \infty$, in order to account for a collection of good subsamples.

4.2. Informed Sub-Sampling MCMC: the methodology

Informed Sub-Sampling MCMC is a scalable adaptation of the Metropolis-Hastings algorithm (Metropolis et al., 1953), designed for situations when N is prohibitively large to perform inference on the posterior π in a reasonable time frame. ISS-MCMC relies on a Markov chain whose transition kernel has a bounded computational complexity, which can be controlled through the parameter n . We first recall how the Metropolis-Hastings algorithm produces a π -reversible Markov chain $\{\theta_i, i \in \mathbb{N}\}$, for any distribution π known up to a normalizing constant.

4.2.1. Metropolis-Hastings

Let Q be a transition kernel on $(\Theta, \mathcal{B}(\Theta))$ and assume that the Metropolis-Hastings Markov chain is at state θ_i . A transition $\theta_i \rightarrow \theta_{i+1}$ consists in the two following step:

- (a) propose a new parameter $\theta \sim Q(\theta_i, \cdot)$
- (b) set the next state of the Markov chain as $\theta_{i+1} = \theta$ with probability

$$a(\theta_i, \theta) = 1 \wedge \alpha(\theta_i, \theta), \quad \alpha(\theta_i, \theta) = \frac{\pi(\theta | Y)Q(\theta, \theta_i)}{\pi(\theta_i | Y)Q(\theta_i, \theta)} \quad (17)$$

and as $\theta_{i+1} = \theta_i$ with probability $1 - a(\theta_i, \theta)$.

Algorithm 1 details how to simulate a Metropolis-Hastings Markov chain $\{\theta_i, i \in \mathbb{N}\}$.

Algorithm 1 Metropolis-Hastings algorithm

```
1: Input: initial state  $\theta_0$  and posterior evaluation  $\pi(\theta_0 | Y)$ 
2: for  $i = 1, 2, \dots$  do
3:   propose a new parameter  $\theta \sim Q(\theta_{i-1}; \cdot)$  and draw  $I \sim \text{unif}(0, 1)$ 
4:   compute  $\pi(\theta | Y)$  and  $a = a(\theta_{i-1}, \theta)$  defined in (17)
5:   if  $I \leq a$  then
6:     set  $\theta_i = \theta$ 
7:   else
8:     set  $\theta_i = \theta_{i-1}$ 
9:   end if
10: end for
11: return: the Markov chain  $\{\theta_i, i \in \mathbb{N}\}$ 
```

4.2.2. Informed Sub-Sampling MCMC

To avoid any confusion, we denote by $\{\tilde{\theta}_i, i \in \mathbb{N}\}$ the sequence of parameters generated by the Informed Sub-Sampling Markov chain, by contrast to the Markov chain $\{\theta_i, i \in \mathbb{N}\}$ produced by the Metropolis-Hastings algorithm (Alg. 1). The pool of *good* subsamples used in the Informed Sub-Sampling inference is treated as a sequence of missing data U_1, U_2, \dots and is thus simulated by our algorithm. More precisely, IIS-MCMC produces a Markov chain $\{(\tilde{\theta}_i, U_i), i \in \mathbb{N}\}$ on the extended space $\Theta \times \mathcal{U}_n$. Inspired by the analysis of Section 3, the sequence of subsamples $\{U_i, i \in \mathbb{N}\}$ is randomly updated in a way that favours those subsets whose summary statistics vector is close to that of the full dataset. Let R be a symmetric transition kernel on $(\mathcal{U}_n, \mathcal{U}_n)$, a transition $(\tilde{\theta}_i, U_i) \rightarrow (\tilde{\theta}_{i+1}, U_{i+1})$ consists in the two following steps:

- i- (a) propose a new subset variable $U \sim R(U_i, \cdot)$
- (b) set $U_{i+1} = U$ with probability

$$b(U_i, U) = 1 \wedge \beta(U_i, U), \quad \beta(U_i, U) = \exp\{\epsilon(\|\Delta_n(U_i)\|^2 - \|\Delta_n(U)\|^2)\} \quad (18)$$

and $U_{i+1} = U_i$ with probability $1 - b(U_i, U)$.

- ii- (a) propose a new parameter $\tilde{\theta} \sim Q(\tilde{\theta}_i, \cdot)$
- (b) set $\tilde{\theta}_{i+1} = \tilde{\theta}$ with probability

$$\tilde{a}(\tilde{\theta}_i, \tilde{\theta}) = 1 \wedge \tilde{\alpha}(\tilde{\theta}_i, \tilde{\theta} | U_{i+1}), \quad \tilde{\alpha}(\tilde{\theta}_i, \tilde{\theta} | U_{i+1}) = \frac{\tilde{\pi}_n(\tilde{\theta} | Y_{U_{i+1}})Q(\tilde{\theta}, \tilde{\theta}_i)}{\tilde{\pi}_n(\tilde{\theta}_i | Y_{U_{i+1}})Q(\tilde{\theta}_i, \tilde{\theta})} \quad (19)$$

and as $\tilde{\theta}_{i+1} = \tilde{\theta}_i$ with probability $1 - \tilde{a}(\tilde{\theta}_i, \tilde{\theta} | U_{i+1})$.

Algorithm 2 details how to simulate an Informed Sub-Sampling Markov chain. Note that at step 11, if $U_i = U_{i-1}$, the quantity $\tilde{\pi}_n(\theta_{i-1} | U_i)$ has already been calculated at the previous iteration.

Algorithm 2 Informed Sub-Sampling MCMC algorithm

```
1: Input: initial state  $(\tilde{\theta}_0, U_0)$  and summary statistics  $S_0 = \bar{S}(Y_{U_0})$ ,  $S^* = \bar{S}(Y)$ 
2: for  $i = 1, 2, \dots$  do
3:   propose a new subset  $U \sim R(U_{i-1}, \cdot)$  and draw  $J \sim \text{unif}(0, 1)$ ,
4:   compute  $S = \bar{S}(Y_U)$  and  $b = b(U_{i-1}, U)$  defined in (18)
5:   if  $J \leq b$  then
6:     set  $U_i = U$  and  $S_i = S$ 
7:   else
8:     set  $U_i = U_{i-1}$  and  $S_i = S_{i-1}$ 
9:   end if
10:  propose a new parameter  $\tilde{\theta} \sim Q(\tilde{\theta}_{i-1}; \cdot)$  and draw  $I \sim \text{unif}(0, 1)$ 
11:  compute  $\tilde{\pi}_n(\tilde{\theta}_{i-1} | Y_{U_i})$ ,  $\tilde{\pi}_n(\tilde{\theta} | Y_{U_i})$  and  $\tilde{a} = \tilde{a}(\tilde{\theta}_{i-1}, \tilde{\theta} | U_i)$  defined in (19)
12:  if  $I \leq \tilde{a}$  then
13:    set  $\tilde{\theta}_i = \tilde{\theta}$ 
14:  else
15:    set  $\tilde{\theta}_i = \tilde{\theta}_{i-1}$ 
16:  end if
17: end for
18: return: the Markov chain  $\{(\tilde{\theta}_i, U_i), i \in \mathbb{N}\}$ 
```

4.3. Connection with noisy ABC

Approximate Bayesian Computation (ABC) is a class of statistical methods, initiated in Pritchard et al. (1999), that allows to infer π in situations where the likelihood f is intractable but forward simulation of pseudo data $Z \sim f(\cdot | \theta)$ is doable. More precisely, the algorithm that consists in (i) $\vartheta \sim p$, (ii) $Z \sim f(\cdot | \vartheta)$ and (iii) set $\theta = \vartheta$ only if $\{Z = Y\}$, does produce a sample θ whose distribution is $\pi(\cdot | Y)$. Regarding the situation $N \rightarrow \infty$ as a source of intractability, one could attempt to borrow from ABC to sample from π . However, since $N \rightarrow \infty$, sampling from the likelihood model is impossible and a natural idea is to replace step (ii) by drawing subsamples Y_U ($U \in \mathcal{U}_n$), leading to what we refer as Informed Sub-Sampling, as opposed to Informed Sub-Sampling MCMC described in the previous Subsection. Obviously, the event $\{Y_U = Y\}$ is impossible except in the trivial situation where $N = n$. Overcoming situations where $\{Y = Z\}$ is impossible or very unlikely has already been addressed in the ABC literature (see Fearnhead and Prangle (2012) and Wilkinson (2013)), leading to noisy ABC algorithms. In particular, step (iii) is replaced by a step that sets $\theta = \vartheta$ with probability $\propto \exp\{-\epsilon \|S(Z) - S(Y)\|^2\}$ where S is a vector of summary statistics and $\epsilon > 0$ a tolerance parameter. We build on this analogy to propose a noisy Informed Sub-Sampling algorithm, see Table 2 for more details.

The Noisy ABC algorithm replaces direct inference of π by the following surrogate distribution

$$\hat{\pi}_{\text{ABC}}(d\theta | Y) : \propto p(d\theta) \hat{f}_{\text{ABC}}(Y | \theta) = p(d\theta) \int f(dZ | \theta) \exp\{-\epsilon \|S(Z) - S(Y)\|^2\}, \quad (20)$$

where the exact likelihood is replaced by \hat{f}_{ABC} . Similarly, the approximation of π stem-

step	ABC		Informed Sub-Sampling	
(i)	$\vartheta \sim p$		-	
(ii)	$Z \sim f(\cdot \vartheta)$		$Z = Y_U, U \sim \text{unif}(U)$	
(iii)	exact	noisy	exact	noisy
	if $Z = Y$, set $\theta = \vartheta$	with proba. \propto $e^{-\epsilon \ S(Y) - S(Z)\ ^2}$ set $\theta = \vartheta$	if $Z = Y$ draw $\theta \sim \pi(U)$	with proba. \propto $e^{-\epsilon \ S(Y) - (N/n)S(Y_U)\ ^2}$ draw $\theta \sim \pi(U)$

Table 2: Comparison between ABC and Informed Sub-Sampling, an adaptation of ABC designed for situations where $N \gg 1$ and likelihood simulation is not possible. The exact algorithms provide samples from π while the noisy algorithms sample from approximation of π given in (20) and (21).

ming from Informed Sub-Sampling is:

$$\hat{\pi}_n(d\theta | Y) : \propto p(d\theta) \hat{f}(Y | \theta) = p(d\theta) \sum_{U \in \mathcal{U}_n} f^{(N/n)}(Y_U | \theta) \exp\{-\epsilon \|(N/n)S(Y_U) - S(Y)\|^2\}. \quad (21)$$

This analogy shows that there is a connection between the ABC and the Informed Sub-Sampling in the way they approximate π , see (20) and (21). However, since sampling from $\nu_{n,\epsilon}$ and $\pi_n(U)$ are not feasible, this approach cannot be considered, hence motivating the use of Markov chains instead *i.e.* Informed Sub-Sampling MCMC. Moreover, quantifying the approximation of π by $\hat{\pi}_n$ (21) is technically challenging while resorting to the Informed Sub-Sampling Markov chain allows to use the Noisy MCMC framework developed in Alquier et al. (2016) to quantify this approximation. This is the purpose of the following Section.

5. Theoretical Analysis of Informed Sub-Sampling MCMC

5.1. Analysis of the Informed Sub-Sampling MCMC algorithm

By construction, IIS-MCMC samples a Markov chain on an extended state space $\{(\tilde{\theta}_i, U_i), i \in \mathbb{N}\}$ but the only useful outcome of the algorithm for inferring π is the marginal chain $\{\tilde{\theta}_i, i \in \mathbb{N}\}$. In this section, we study the distribution of the marginal chain and denote by $\tilde{\pi}_i$ the distribution of the random variable $\tilde{\theta}_i$. Note that $\{\tilde{\theta}_i, i \in \mathbb{N}\}$ is identical to the Metropolis-Hastings chain $\{\theta_i, i \in \mathbb{N}\}$, up to replacing α by $\tilde{\alpha}$ in the accept/reject step. This change, from which the computational gain of our method originates, implies that π is not the stationary distribution of $\{\tilde{\theta}_i, i \in \mathbb{N}\}$. Interest lies in quantifying the distance between $\tilde{\pi}_i$ and π . We first recall the definition of the total variation distance, for two distributions with density function π and $\tilde{\pi}_i$ respectively w.r.t. the same common dominating measure,

$$\|\pi - \tilde{\pi}_i\|_{\text{TV}} = (1/2) \int |\pi(\theta) - \tilde{\pi}_i(\theta)| d\theta.$$

We cast the analysis of the marginal chain $\{\tilde{\theta}_i, i \in \mathbb{N}\}$ in the noisy Markov chain Monte Carlo framework developed in Alquier et al. (2016). Our main result is the following proposition:

Proposition 4. *Let K be the transition kernel of the marginal Metropolis-Hastings Markov chain that have acceptance ratio α . Assume K is uniformly ergodic, i.e there exists two constants $C < \infty$ and $\varrho < 1$ such that for all $i \in \mathbb{N}$*

$$\sup_{\theta_0 \in \Theta} \|\pi - K(\theta_0, \cdot)^i\|_{\text{TV}} \leq C\varrho^i.$$

Let \tilde{K} be the transition kernel of the marginal Informed Sub-Sampling Markov chain. Define

$$D_1(U, \theta) := \mathbb{E}\{\alpha(\theta, \theta')|\phi_U(\theta) - \phi_U(\theta')|\} = \int Q(\theta, d\theta')\alpha(\theta, \theta')|\phi_U(\theta) - \phi_U(\theta')|, \quad (22)$$

where for all $(\theta, U) \in (\Theta \times \mathbf{U}_n)$, we have set $\phi_U(\theta) = f(Y_U | \theta)^{N/n} / f(Y | \theta)$ and

$$D_2(\theta) := \mathbb{E}\left\{\frac{f(Y | \theta)}{f(Y_U | \theta)^{N/n}}\right\} = \sum_{U \in \mathbf{U}_n} \nu_{n, \epsilon}(U) \frac{f(Y | \theta)}{f(Y_U | \theta)^{N/n}}. \quad (23)$$

Then,

- we have:

$$\lim_{i \rightarrow \infty} \|\pi - \tilde{\pi}_i\|_{\text{TV}} \leq \kappa \sup_{\theta \in \Theta} \sup_{i \in \mathbb{N}} \mathbb{E}_i \left\{ \frac{f(Y | \theta)}{f(Y_U | \theta)^{N/n}} \right\} \sup_{U \in \mathbf{U}_n} D_1(U, \theta), \quad (24)$$

where $\kappa = \lambda + C\varrho^\lambda / (1 - \varrho)$ is a constant that depends only on the Metropolis-Hastings chain and $\lambda = \lceil \log_\varrho(1/C) \rceil$ and \mathbb{E}_i is the expectation with respect to p_i defined as the distribution of U_i .

- if in addition the IIS-MCMC initial distribution for the subset variable U_0 is set to $p_0 = \nu_{n, \epsilon}$, then we have:

$$\lim_{i \rightarrow \infty} \|\pi - \tilde{\pi}_i\|_{\text{TV}} \leq \kappa \sup_{\theta \in \Theta} D_2(\theta) \sup_{U \in \mathbf{U}_n} D_1(U, \theta). \quad (25)$$

The proof of Proposition 4 is postponed to [Appendix A.4](#).

Since for any two measures (μ, μ') , $\|\mu - \mu'\|_{\text{TV}} \leq 1$, the upper bounds of Proposition 4 are only informative if there are smaller than 1. Those bounds are a product of two expectations. We show

- in section 5.2, how the choice of proposal kernel Q allows to take D_1 arbitrarily close to 0.
- in section 5.3, under the assumption $p_0 = \nu_{n, \epsilon}$, that conditions on the parameters of the Informed Sub-Sampling MCMC algorithm, namely the subset size n , the bandwidth ϵ and the summary statistics S , control (i.e bound) D_2 .

The central piece of our work has been to show that injecting a perturbation consisting of replacing Y by Y_U and f by $f^{N/n}$ into the acceptance ratio of an uniformly ergodic Metropolis-Hastings Markov chain targeting π is asymptotically absorbed when the choice of subsamples is guided by $\nu_{n, \epsilon}$ and the summary statistics satisfy a condition detailed in section 5.3.

5.2. On the Proposal Kernel

Assuming a Gaussian random walk proposal with covariance matrix $\Sigma^T \Sigma$, D_1 writes

$$D_1(U, \theta) = \int \Phi_d(d\zeta) \frac{\pi(\theta + \Sigma\zeta)}{\pi(\theta)} |\phi_U(\theta) - \phi_U(\theta + \Sigma\zeta)|, \quad (26)$$

where Φ_d is the standard Gaussian distribution in dimension $d = \dim(\Theta)$. When $N \gg 1$, the Bernstein-von Mises theorem states that, under conditions on the likelihood function, the posterior distribution can be approximated by a Gaussian with mean set as the maximum likelihood estimator θ^* and covariance $I(\theta_0)^{-1}/N$ where I is the Fisher information matrix and $\theta_0 \in \Theta$ some parameter. Since ISS-MCMC aims at sampling from an approximation of π , setting $\Sigma = (1/\sqrt{N})M$ where $M^T M$ is an approximation of $I(\theta_0)^{-1}$ is a reasonable choice. Proposition 5 shows that D_1 can, in this scenario, be arbitrarily brought down close to 0.

Proposition 5. *Under the assumption that the proposal kernel Q is a Gaussian Random Walk with covariance matrix $\Sigma = (1/\sqrt{N})M$ then*

$$D_1(U, \theta) \leq \frac{\|\nabla_\theta \phi_U(\theta)\|}{\sqrt{N}} \left\{ \sqrt{\frac{2}{\pi}} \|M\|_1 + \frac{\|M\|_2^2 \|\nabla_\theta \log \pi(\theta)\|}{\sqrt{N}} \right\} \\ + \frac{d}{2N} \|M^T \nabla_\theta^2 \phi_U(\theta) M\| + \mathbb{E}\{R(\|M\zeta\|/\sqrt{N})\}, \quad (27)$$

where $R(x) =_{x \rightarrow 0} o(x)$ and for any square matrix M of dimension \mathbb{R}^d , we have set $\|M\|_1 := \sum_{1 \leq i, j \leq d} |M_{i,j}|$, $\|M\|_2 := \{\sum_{1 \leq i, j \leq d} M_{i,j}^2\}^{1/2}$ and $\|M\|$ is the operator norm of M .

The proof is postponed to [Appendix A.6](#). Under regularity assumptions on the likelihood model, the gradient of $\log \pi$ and ϕ_U and the Hessian of ϕ_U are bounded and the upper bound of $D_1(U, \theta)$ can be brought down arbitrarily to 0, uniformly in (U, θ) , through M when $N \gg 1$.

5.3. On Summary Statistics

In this paragraph, we assume that $p_0 = \nu_{n,\epsilon}$ i.e the auxiliary chain $\{U_k, k \in \mathbb{N}\}$ starts at stationarity and (25) holds. Proposition 5, stated at the previous subsection, allows to assess the quality of ISS-MCMC only in the case where the expectation

$$D_2(\theta) = \mathbb{E} \left\{ \frac{f(Y|\theta)}{f(Y_U|\theta)^{N/n}} \right\} = \sum_{U \in \mathcal{U}_n} \nu_{n,\epsilon}(U) \frac{f(Y|\theta)}{f(Y_U|\theta)^{N/n}} \quad (28)$$

is bounded. Proposition 3 shows that in case of an exponential model where a sufficient statistics mapping \bar{S} exists, our definition of $\nu_{n,\epsilon}$:

$$\nu_{n,\epsilon}(U) \propto \exp\{-\epsilon \|\Delta_n(U)\|^2\}, \quad \Delta_n(U) = N(\bar{S}(Y) - \bar{S}(Y_U)), \quad (29)$$

is sufficient to bound $D_2(\theta)$ for all $\theta \in \Theta$. Here, \bar{S} is a non-extensive mapping $\bar{S} : \mathcal{Y}^p \rightarrow \mathbb{R}^s$ from datasets of any size $p > 0$ to the sufficient statistics space. In the context of Proposition 3, for any $Y \in \mathcal{Y}^p$, $\bar{S}(Y) = (1/p) \sum_{i=1}^p \bar{S}(Y_i)$. Following this result and extending

the definition of $\nu_{n,\epsilon}$ beyond the case of exponential models by replacing sufficient statistics by summary, we propose a verifiable condition on the summary statistics ensuring that $D_2(\theta)$ is bounded.

In Eq. (28), the likelihood of each subsample is raised at the power N/n (i.e. typically several orders of magnitude) and therefore subsamples unlikely under $f(\cdot|\theta)$ will contribute to make $D_2(\theta)$ very large. Ideally the choice of S would guarantee that subsamples Y_U having a very small likelihood $f(Y_U|\theta)$ are assigned to a weight $\nu_{n,\epsilon}(U) \approx 0$ to limit their contribution. In other words, S should be specified in a way that prevents $f(Y_U|\theta)$ to go to 0 at a rate faster than $\nu_{n,\epsilon}(U)$. This leads to the following condition:

Condition 1. *There exists a summary statistics mapping S , a constant $\gamma > 0$, such that for all $(\theta, U) \in \Theta \times \mathcal{U}_n$*

$$\log f(Y|\theta) - (N/n) \log f(Y_U|\theta) \leq \gamma N \|\bar{S}(Y) - \bar{S}(Y_U)\|.$$

Indeed, if such a condition holds,

$$\begin{aligned} D_2(\theta) &= \sum_{U \in A_n(\theta)} \nu_{n,\epsilon}(U) \frac{f(Y|\theta)}{f(Y_U|\theta)^{N/n}} + \sum_{U \in \mathcal{U}_n \setminus A_n(\theta)} \nu_{n,\epsilon}(U) \frac{f(Y|\theta)}{f(Y_U|\theta)^{N/n}} \\ &\leq \nu_{n,\epsilon}(A_n(\theta)) + \sum_{U \in \mathcal{U}_n \setminus A_n(\theta)} \exp\{-\epsilon \|\Delta_n(U)\|^2 + \gamma \|\Delta_n(U)\| - \log Z(\epsilon)\}, \end{aligned}$$

where, similarly to Proposition 3, we have defined for all $\theta \in \Theta$, $A_n(\theta) = \{U \in \mathcal{U}_n, f(Y|\theta) < f(Y_U|\theta)^{N/n}\}$ and $Z(\epsilon) = \sum_{U \in \mathcal{U}_n} \exp\{-\epsilon \|\Delta_n(U)\|^2\}$. Clearly, if ϵ has the same order of magnitude as γ , each term of the sum remains bounded when $\|\Delta_n(U)\| \rightarrow \infty$. Conversely, setting $\epsilon = 0$ is equivalent to choosing $\nu_{n,\epsilon}$ as the uniform distribution on \mathcal{U}_n and may not allow to bound $D_2(\theta)$, see Figure 3 related to the probit example.

Potential summary statistics can be validated by checking that they satisfy Condition 1. This validation can be performed graphically, by repeating the following operations for a number of parameters $\theta_k \sim_{\text{i.i.d}} p$:

- (i) draw subsets U_1, U_2, \dots uniformly at random in \mathcal{U}_n ,
- (ii) plot the points with coordinates

$$(x_{k,i}, y_{k,i}) = (\|\Delta_n(U_i)\|, \log f(Y|\theta_k) - (N/n) \log f(Y_{U_i}|\theta_k)).$$

The statistics are validated if there exists $\gamma < \infty$ such that the points $(x_{k,i}, y_{k,i})$ satisfy $|y_{k,i}/x_{k,i}| \leq \gamma$, see for example Figure 4.

In situations where the maximum likelihood estimator $\theta^*(Y_{1:n})$ is easy and quick to evaluate numerically, we recommend setting $\bar{S}(Y_{1:n}) = \theta^*(Y_{1:n})$. In the case of independent observations of a well-specified model, setting the summary statistics as the maximum likelihood estimate is justified by the following Proposition which implies that Condition 1 holds, asymptotically.

Proposition 6. *We assume that the whole dataset comprises $N = pn$ independent observations and there exists some $\theta_0 \in \Theta$ such that $Y_i \sim f(\cdot|\theta_0)$. Let θ^* be the MLE of Y_1, \dots, Y_N and θ_U^* be the MLE of the subsample Y_U ($U \in \mathcal{U}_n$). Then, there exists a*

constant β , a metric $\|\cdot\|_{\theta_0}$ on Θ and a non-decreasing subsequence $\{\sigma_n\}_{n \in \mathbb{N}}$, ($\sigma_n \in \mathbb{N}$) such that for all $U \subset \{1, 2, \dots, \rho\sigma_n\}$ with $|U| = \sigma_n$, we have for p -almost all θ in an neighborhood of θ_0 :

$$\log f(Y_{1:\rho\sigma_n} | \theta) - \rho \log f(Y_U | \theta) \leq H_n(Y, \theta) + \beta + \frac{\rho n}{2} \|\theta_U^* - \theta^*\|_{\theta_0}, \quad (30)$$

where

$$\text{plim}_{n \rightarrow \infty} H_n(Y, \theta) \stackrel{\mathbb{P}_{\theta_0}}{=} 0.$$

The proof is postponed to [Appendix B](#) and follows from a careful application of a Bernstein-von Mises theorem. Note that extension of [Proposition 6](#) to cases where the observations are not independent may exist provided that a Bernstein-von Mises theorem holds for the model at hand, which is the case for dependent observations if the likelihood model satisfies local asymptotic normality conditions [Le Cam \(1953, 1986\)](#).

We remark that [Proposition 6](#) is in line with the results regarding optimal summary statistics for ABC established in [Fearnhead and Prangle \(2012\)](#). The authors show that the quadratic error loss between the ABC estimate based on $\hat{\pi}_{\text{ABC}}$ ([Eq. 20](#)) and the true parameter is minimized when setting the summary statistics as the posterior mean, a choice which asymptotically coincides with the maximum likelihood estimator.

6. Illustrations

We evaluate the efficiency of ISS-MCMC on three different applications: inferring a time series observed at $N = 10^6$ contiguous time steps, a logistic regression with $N = 10^6$ observations and a Gaussian binary classification problem based on $N = 10^7$ data.

6.1. Implementation details of Informed Sub-Sampling MCMC

Before illustrating the ISS-MCMC algorithm on the different examples, we address a few technical implementation details.

- On the subset size n : this parameter is essentially related to the computational budget available to the user. In the following examples we have used $n \propto N^{1/2}$ which achieves a substantial computational gain at a price of a negligible asymptotic bias.
- On the sufficient statistics S : to reduce the bias resulting from the Metropolis-Hastings approximation, S should be constructed so that [Condition 1](#) holds. If the maximum likelihood estimator $\theta^*(Y)$ is quick to compute then [Proposition 6](#) suggests that setting $S(Y) = \theta^*(Y)$ will theoretically satisfy [Condition 1](#). Other sufficient statistics mapping can be used, typically those arising in the Approximate Bayesian Computation literature. In any case, we recommend checking [Condition 1](#) graphically (see [Section 5](#)).
- On the bandwidth parameter ϵ : the theory shows that when $\epsilon \approx \gamma N^2$, the asymptotic bias is controlled (γ is the constant in [Condition 1](#)). In practice, this may prove to be too large and could potentially cause the algorithm to get stuck on a very small number of subsets. To avoid such a situation, we suggest monitoring the refresh rate of subsamples that should occur with probability of at least 1%.

- On the initial subset U_0 : In theory, one would run a preliminary Markov chain $\{U_1^{(0)}, \dots, U_L^{(0)}\}$ (for some $L > 0$) targeting $\nu_{n,\epsilon}$, and set $U_0 = U_L^{(0)}$ in order for the results of Section 5.3 to hold. In practice, a more efficient approach is to use a simulating annealing Metropolis-Hastings algorithm, see [Geyer and Thompson \(1995\)](#). It introduces a sequence of tempered distributions $\nu_k := \nu_{n,\epsilon_k}$, such that $\epsilon_k = t_k \epsilon$ ($k \in \{1, \dots, L\}$) where $t_1 = 0$ and $t_L = 1$. The transition kernel of the k -th iteration of the preliminary Markov chain is designed to be ν_k invariant. This technique facilitates sampling from a proxy of $\nu_{n,\epsilon}$ in a relative short time period as the successive tempered distributions help identifying those subsamples belonging to the high probability sets of $\nu_{n,\epsilon}$.

6.2. Inference of an AR(2) model

Example 3. An autoregressive time series of order 2 AR(2) $\{Y_k, k \leq N\}$ is defined recursively by:

$$\begin{cases} (Y_0, Y_1) \sim \mu := \mathcal{N}_2(\mathbf{0}_2, \theta_3^2 Id_2) \\ Y_n = \theta_1 Y_{n-1} + \theta_2 Y_{n-2} + Z_n, \quad Z_n \sim \mathcal{N}(0, \theta_3^2), \quad \forall n \geq 2, \end{cases} \quad (31)$$

where $\theta \in \Theta \subset \mathbb{R}^3$. The likelihood of an observed time series for this model writes

$$f(Y_{0:N} | \theta) = \mu(Y_{0:1}) \prod_{k=2}^N g_k(Y_k | Y_{0:k-1}, \theta), \quad (32)$$

such that for all $k \geq 1$,

$$g(Y_k | Y_{0:k-1}, \theta) = \Phi_1(Y_k; \theta_1 Y_{k-1} + \theta_2 Y_{k-2}, \theta_3^2) \quad (33)$$

where $x \rightarrow \Phi_1(x; m, v)$ is the pdf of the univariate Gaussian distribution with mean m and variance v .

This model has been used in [Chib and Greenberg \(1995, Section 7.2\)](#) to showcase the Metropolis-Hastings (M-H) algorithm. We follow the same setup here and in particular we use the same true parameter $\theta^* = (1, -0.5, 1)$, same prior distribution and proposal kernel Q ; see [Chib and Greenberg \(1995\)](#) for more details. We sampled a time series $\{Y_k, k \leq N\}$ according to (31), with $N = 10^6$ under θ^* . Of course, in such a setup, M-H is prohibitively slow to be used in practice to sample from π as it involves evaluating the likelihood of the whole time series at each iteration. We nevertheless use M-H to obtain a ground truth of π .

For simplicity, we restrict the set of subsamples to n contiguous observations:

$$\{Y_{0:n-1}, Y_{1:n}, \dots, Y_{N-n+1:N}\}.$$

This induces a set of subset $\bar{U}_n \subset U_n$ defined such that a subset $U \in \bar{U}_n$ is identified with its starting index, *i.e* for all $i \leq |\bar{U}_n|$, $\bar{U}_n \ni U_i := \{i, i+1, \dots, i+n-1\}$. Indeed, using such subsamples yields a tractable likelihood (32) as otherwise, missing variables need to be integrated out, hence losing the simplicity of our approach.

With some abuse of notations, the proposal kernel R can be written as a transition kernel on the alphabet $\{0, \dots, N - n + 1\}$. It is defined in this example as:

$$R(i; j) = \mathbf{1}_{i \neq j} \left\{ \omega \frac{\exp(-\lambda|j - i|)}{\sum_{j \leq |\bar{U}_n|, j \neq i} \exp(-\lambda|j - i|)} + (1 - \omega) \frac{1}{|\bar{U}_n| - 1} \right\}. \quad (34)$$

The rationale is to propose a new subset through a mixture of two distributions: the first gives higher weight to local moves and the latter allows jumps to remote sections of the time series. In this example, we have used $\omega = 0.9$ and $\lambda = 0.1$. We study the efficiency of ISS-MCMC in function of n , ϵ and S .

For any subsample Y_U , $U \in \bar{U}_n$, we have set the summary statistics $\bar{S}(Y_U)$ to the solution of the AR(2) Yule-Walker equations for the dataset Y_U . As shown in Figure 4, this choice of summary statistics satisfies (graphically) Condition 1 with $\gamma \approx 5 \cdot 10^6$. We therefore set $\epsilon = 5.0 \cdot 10^6$ to make sure that D_2 (23) is bounded. Theoretically, Proposition 4 guarantee that the bias is controlled. This is illustrated graphically in Figure 5 where π is compared to $\tilde{\pi}_i$ ($i = 50,000$). We also report the distribution of the Informed Sub-Sampling chain when the $\epsilon = 0$, *i.e* when the subsampling is actually uninformed and all subsamples have the same weight. In the latter case, D_2 is not bounded which explains why the bias on $\lim_{i \rightarrow \infty} \|\pi - \tilde{\pi}_i\|_{\text{TV}}$ is not controlled. Figures 5 and 6 illustrate the distribution of $\{\tilde{\theta}_i, i \in \mathbb{N}\}$ for some runs of ISS-MCMC with $\epsilon = 0$ and $\epsilon = 5.0 \cdot 10^6$. Finally, Figure 9 gives a hint at the computational efficiency of ISS-MCMC. Metropolis-Hastings was compared to IIS-MCMC with $n \in \{1,000; 5,000; 10,000\}$ and $\epsilon \in \{0; 1; 5.0 \cdot 10^6\}$. The performance indicator is defined as the average of the marginals Total Variation distance, *i.e*

$$\text{TV}(t) = \frac{1}{d} \sum_{j=1}^d \|\pi^{(j)} - \tilde{\pi}_t^{(j)}\|_{\text{TV}},$$

where $\pi^{(j)}$ and $\tilde{\pi}_t^{(j)}$ are respectively the true j -th marginal and the j -th marginal of the chain distribution after a runtime of t seconds. The true marginals were estimated from a long Metropolis-Hastings chain, at stationarity. $\tilde{\pi}_t^{(j)}$ was estimated using 500 independent chains starting from the prior. The Matlab function `ksdensity` in default settings was applied to estimate $\|\pi^{(j)} - \tilde{\pi}_t^{(j)}\|_{\text{TV}}$ from the chains samples, hence the variability. On the one hand, when $\epsilon = 0$, there is no informed search for subsamples which makes the algorithm much faster than the other setups but yields a significant bias (larger than 0.5). On the other hand, setting $\epsilon > 0$ adds to the computational burden but allows to reduce the bias. In fact, for $n = 5,000$ and a computational budget of $t = 1,000$ seconds, the bias of ISS-MCMC is similar to that of Metropolis Hastings but converges 100 times as fast. Finally, note that as expected by the theory, when the unrepresentative subsamples are not penalized enough (*e.g.* by setting $\epsilon = 1$), ISS-MCMC yields a significant bias which hardly improves on uninformed subsampling when $n = 10,000$. Following Condition 1, we see that setting $\epsilon = 5.0 \cdot 10^6$ significantly reduces the bias. Note that setting $\epsilon > 5.0 \cdot 10^6$ could potentially reduce further the bias but may fail the algorithm: indeed, when ϵ is too large the chain $\{U_k, k \in \mathbb{N}\}$ gets easily stuck on a set of the best subsamples (for this choice of summary statistics) and may considerably slow down the convergence of the algorithm.

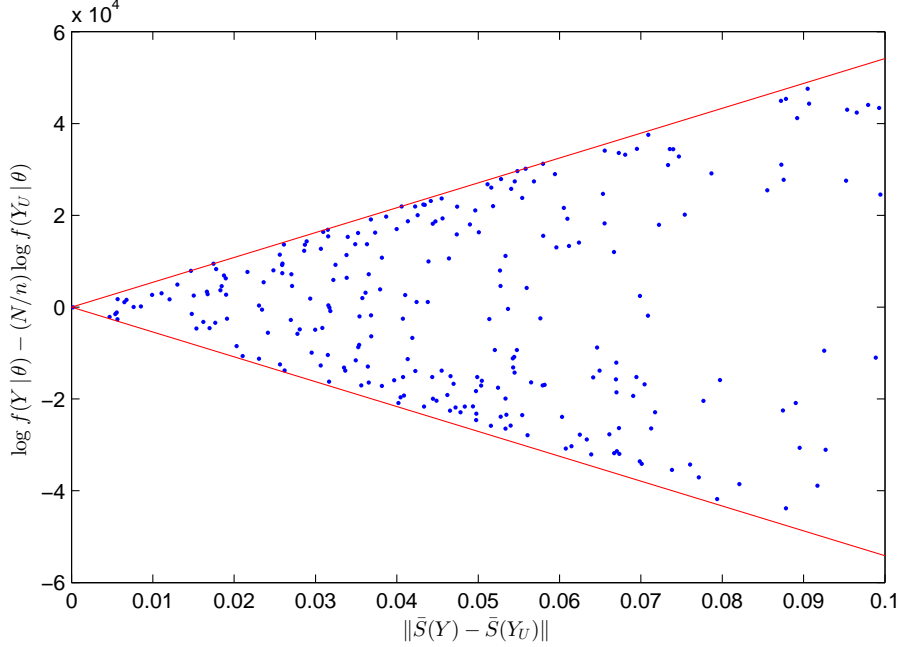


Figure 4: (Example 3: Autoregressive time series) Validation of summary statistics set as the solution of Yule-Walker equation, with $n = 1,000$. This choice of sufficient statistics satisfies Condition 1 with $\gamma \approx 5 \cdot 10^6$. Each dot corresponds to a point $(\log f(Y | \theta) - (N/n) \log f(Y_U | \theta), \|\bar{\Delta}_n(U)\|)$ where θ and U were respectively drawn from the prior p and uniformly at random in \mathbf{U}_n . The red lines allow to estimate the parameter γ of Condition 1.

Other choices of summary statistics can be considered. Since Y is modelled as an autoregressive time series, an option would be to set the summary statistics as the empirical autocorrelation function. Figure 8 shows that it is not a recommended choice. The left panel suggests that Condition 1 does not hold for this type of summary statistics: good subsets yield a large value for $\log f(Y | \theta) - (N/n) \log f(Y_U | \theta)$ and conversely for bad subsets, hence generating a bias (see Eq. (28)). As a consequence, the right panel which shows a clear mismatch between $\pi^{(3)}$ and the Informed Sub-Sampling third marginal is not surprising.

6.3. Logistic Regression Example

Example 4. A d -dimensional logistic regression model is parameterized by a vector $\theta = (\theta_1, \dots, \theta_d) \in \Theta \subset \mathbb{R}^d$. Observations are realizations of the following model:

- simulate covariates $X_i = (X_{i,1}, \dots, X_{i,d}) \sim \mathcal{N}(0, (1/d)^2)$
- simulate Y_i given θ and X_i as

$$Y_i = \begin{cases} 1 & \text{w.p. } 1 / (1 + e^{-\theta X^T}), \\ 0 & \text{otherwise.} \end{cases} \quad (35)$$

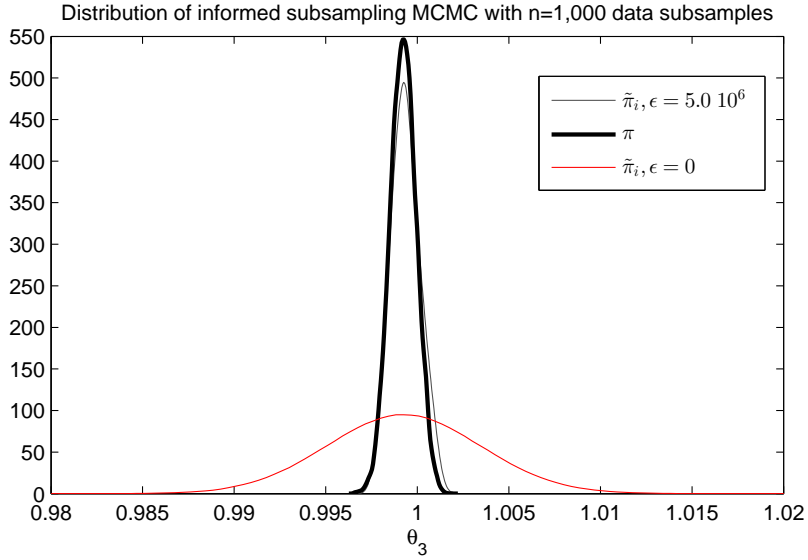


Figure 5: (Example 3: Autoregressive time series) Inference of the noise parameter with ISS-MCMC, using subsets comprising of $n = 1,000$ contiguous time steps of a $N = 10^6$ time-series. The plot represents the distributions $\tilde{\pi}_i$ ($i = 50,000$) of the Informed Sub-Sampling Markov chain for two different values of $\epsilon \in \{0, 5.0 \cdot 10^6\}$. These distributions were obtained from the replication of 1,000 independent chains.

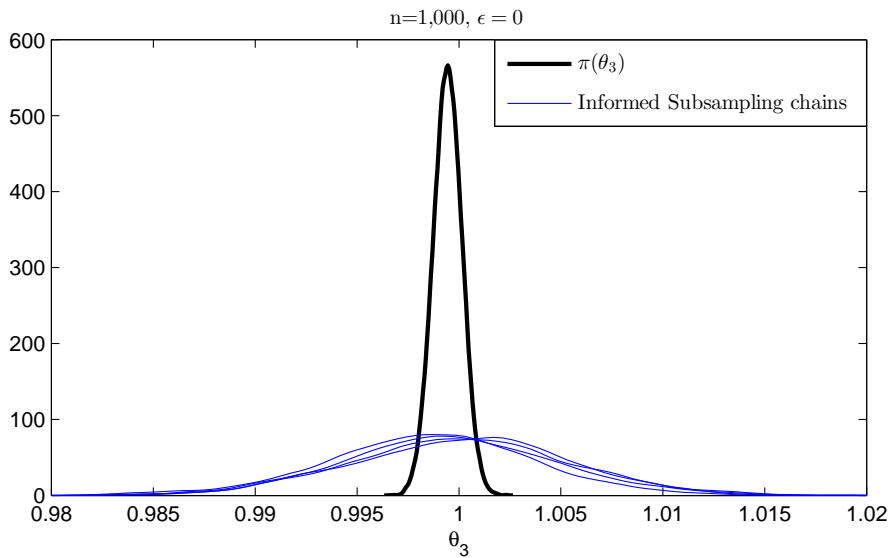


Figure 6: (Example 3: Autoregressive time series) Marginal distribution of θ_3 and distribution of $\{\tilde{\theta}_i, i \in \mathbb{N}\}$ for four independent Informed Sub-Sampling Markov chains with $\epsilon = 0$ and $n = 1,000$.

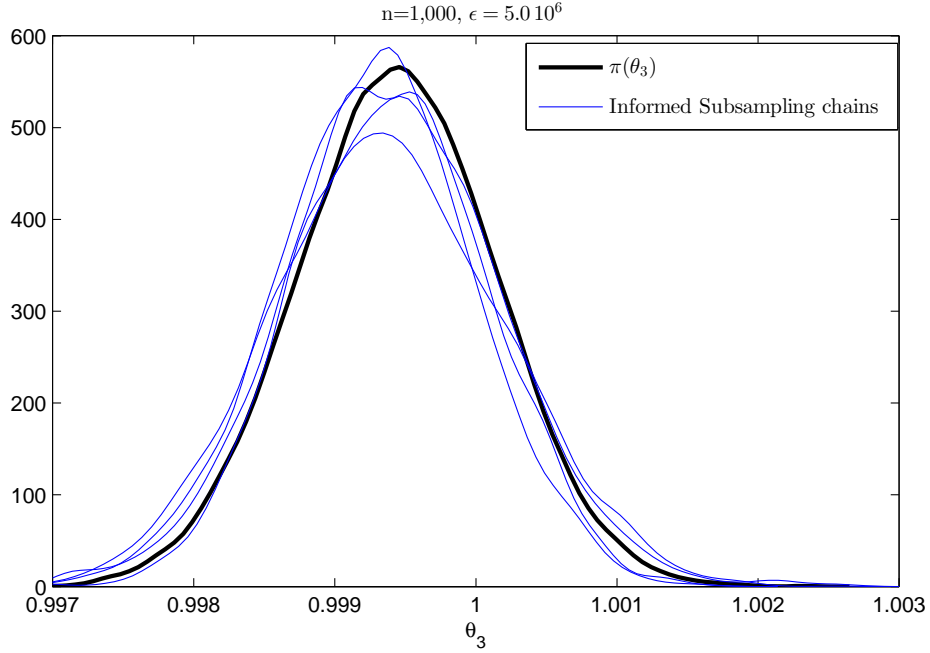


Figure 7: (Example 3: Autoregressive time series) Marginal distribution of θ_3 and distribution of $\{\tilde{\theta}_i, i \in \mathbb{N}\}$ for four independent Informed Sub-Sampling Markov chains with $\epsilon = 5.0 \cdot 10^6$ and $n = 1,000$.

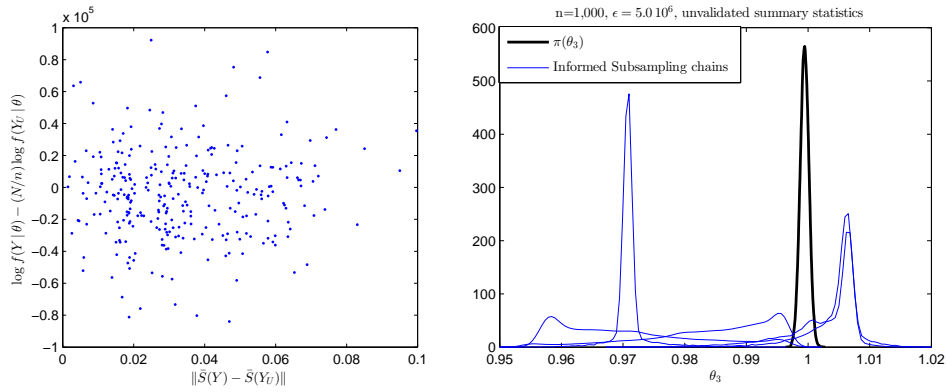


Figure 8: (Example 3: Autoregressive time series) In this case, the summary statistics were defined as the first 5 empirical autocorrelation coefficients. The left panel shows that this is not a recommended choice and the right panel illustrates the distribution of $\{\tilde{\theta}_i, i \in \mathbb{N}\}$ for four independent Informed Sub-Sampling Markov chains ($\epsilon = 5.0 \cdot 10^6$, $n = 1,000$ and this choice of summary statistics), yielding an obvious mismatch.

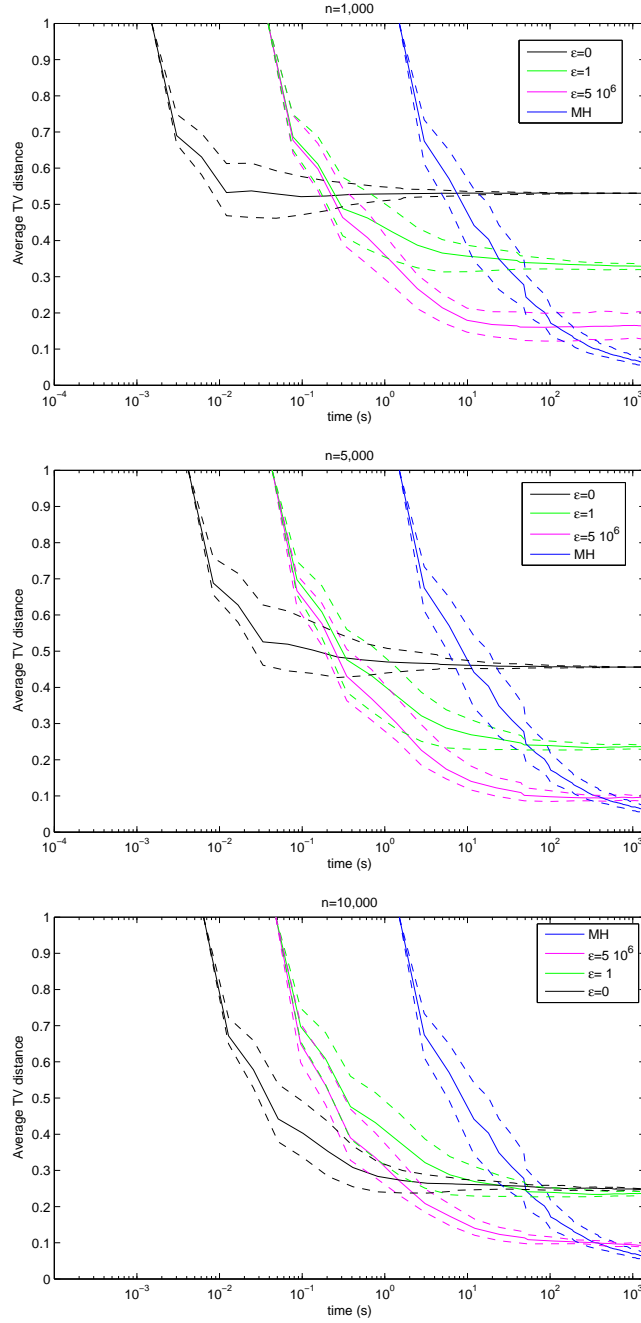


Figure 9: (Example 3: Autoregressive time series) Average Total variation distance over the three marginals between π and $\tilde{\pi}_t$ in function of the simulation time t . The dashed lines represent the first and third quartiles. Scenario $n = 1,000$ (top), $n = 5,000$ (middle) and $n = 10,000$ (bottom) with three different ϵ . Note that in all the three plots, the MH curves are identical and are just reported for comparison purpose.

algorithm	time/iter.(s)	iter. completed	RMSE	var $\{\widehat{\pi(D)}\}$
M-H	10	50	0.1417	0.004
IIS-MCMC, $n = 1,000$	0.05	10,000	0.1016	0.0104
IIS-MCMC, $n = 5,000$	0.08	6,250	0.0351	0.0012
IIS-MCMC, $n = 10,000$	0.13	3,840	0.0267	0.0007

Table 3: (Example 4: Logistic regression) Tradeoff Bias-Variance of the Monte Carlo estimator from Metropolis-Hastings and ISS-MCMC for a fixed computational budget of 500 seconds. Those results were replicated using 100 replications of each algorithm.

We have simulated $N = 10^6$ observations Y_1, Y_2, \dots under the true parameter $\theta^* = (1, 2, -1)$ ($d = 3$).

The summary statistics were set as the maximum likelihood estimator returned by the Matlab routine `glmfit` and were graphically validated, as in Figure 4. The tolerance parameter was consequently set $\epsilon = 5.0 \cdot 10^{-4}$. We study the influence of n on the Informed Sub-Sampling chain marginal distributions in Figure 10. We note that as soon as $n \geq 5,000$, the bias vanishes and that when random subsampling is used (*i.e.* $\epsilon = 0$), the bias is much larger. Of course, Figure 10 only gives information about the marginal distributions. To complement the study, we consider estimating the probability $\pi(D)$ where D is the domain defined as:

$$D \subset \Theta = \{\theta_1 \in (0.98, 1.00); \theta_2 \in (1.98, 2.01); \theta_3 \in (-0.98, -0.95)\},$$

in order to check that the joint distribution π is reasonably inferred. Numerical integration using a long Metropolis-Hastings algorithm, gave the ground truth $\pi(D) = 0.1$. The top panel of Figure 11 illustrates the Monte Carlo estimation of $\pi(D)$ based on $i = 10,000$ iterations of ISS-MCMC implemented with $n \in \{1,000; 5,000; 10,000\}$ and compares it to Metropolis-Hastings. As expected ISS-MCMC has a negligible bias and the variance of the estimator decreases when n increases. Indeed, when n increases, the Informed Sub-Sampling process is less likely to pick irrelevant subsets, which in turns lower the variability of the chain. The Monte-Carlo estimation based on ISS-MCMC with $n = 10,000$ and Metropolis-Hastings are very similar. However, when we normalize the experiment by the CPU time, Metropolis-Hastings is clearly outperformed by ISS-MCMC. The lower panel of Figure 10 assumes that only $t = 500$ seconds of computation are available. All the chains are started from θ^* . Table 3 reports the quantitative details of this experiment. In such a situation, one should clearly opt for the Informed Sub-Sampling approach as the Metropolis-Hastings algorithm only achieves 50 iterations for this amount of computation and as such fails to reach stationarity.

6.4. Binary Classification

Example 5. A training dataset consisting of $N = 10^7$ labeled observations $Y = \{Y_k, k \leq N\}$ from a 2 dimensional Gaussian mixture model is simulated with

$$Y_k | I_k = i \sim \mathcal{N}(\mu_i, \Gamma_i), \quad I_k \sim \text{Bernoulli}(1/2),$$

where $\mu_1 = [\theta_1, 0]$, $\mu_2 = [\theta_2, 0]$, $\Gamma_1 = \text{diag}([\theta_3/2, \theta_3])$ and $\Gamma_2 = \text{diag}([\theta_4/2, \theta_4])$. We define $\theta = (\theta_1, \theta_2)$ with $\theta_i \in \mathbb{R} \times \mathbb{R}_*^+$ for each model $i \in \{1, 2\}$. A prior distribution

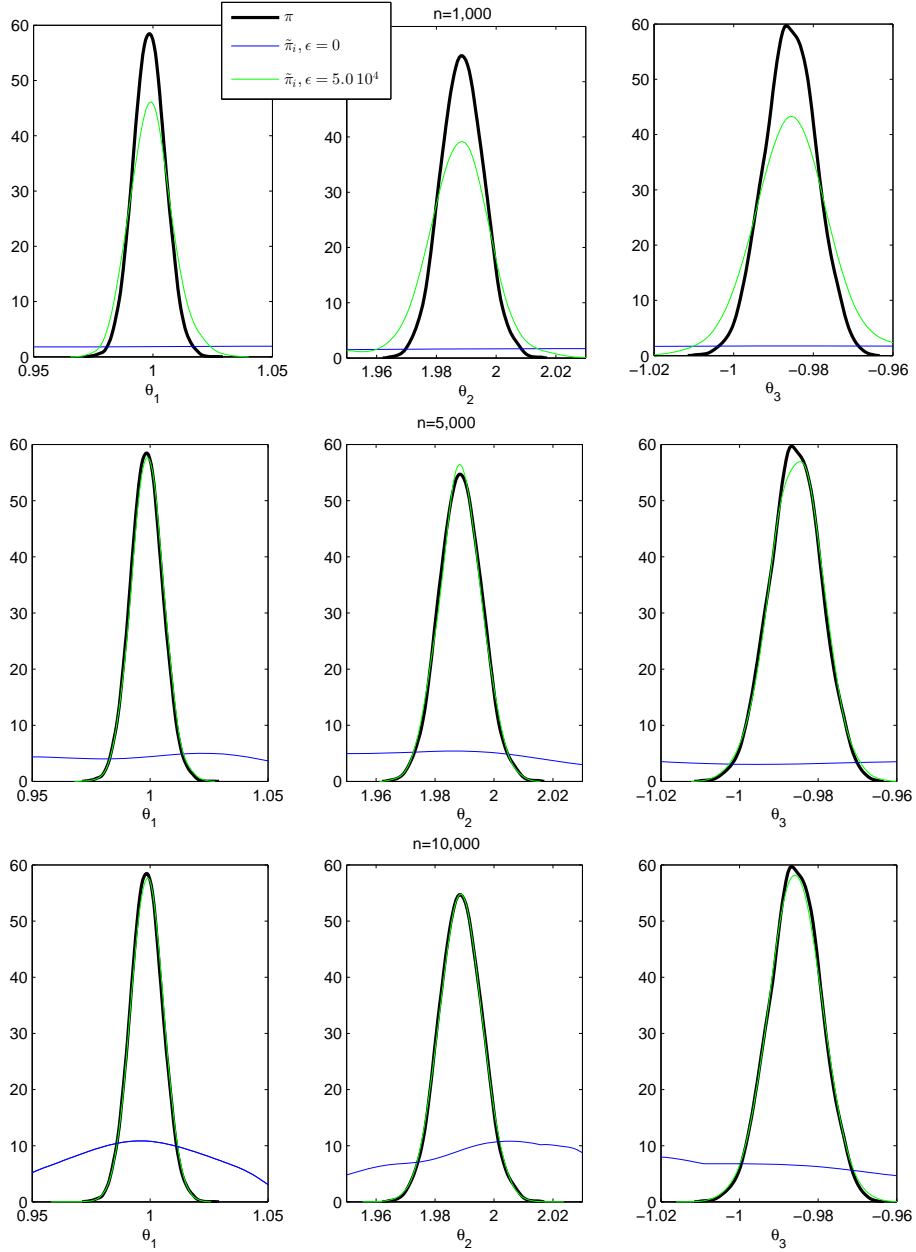


Figure 10: (Example 4: Logistic regression) Stationary marginal distributions of the ISS-MCMC algorithm $\tilde{\pi}_i$, (for $i = 1,000$) with $\epsilon = \{0; 5.0 \cdot 10^{-4}\}$ (respectively in blue and green) and $n \in \{1,000; 5,000; 10,000\}$ and true marginal π (in black). $\tilde{\pi}_i$ was estimated by simulating 1,000 copies of Informed Sub-Sampling chains.

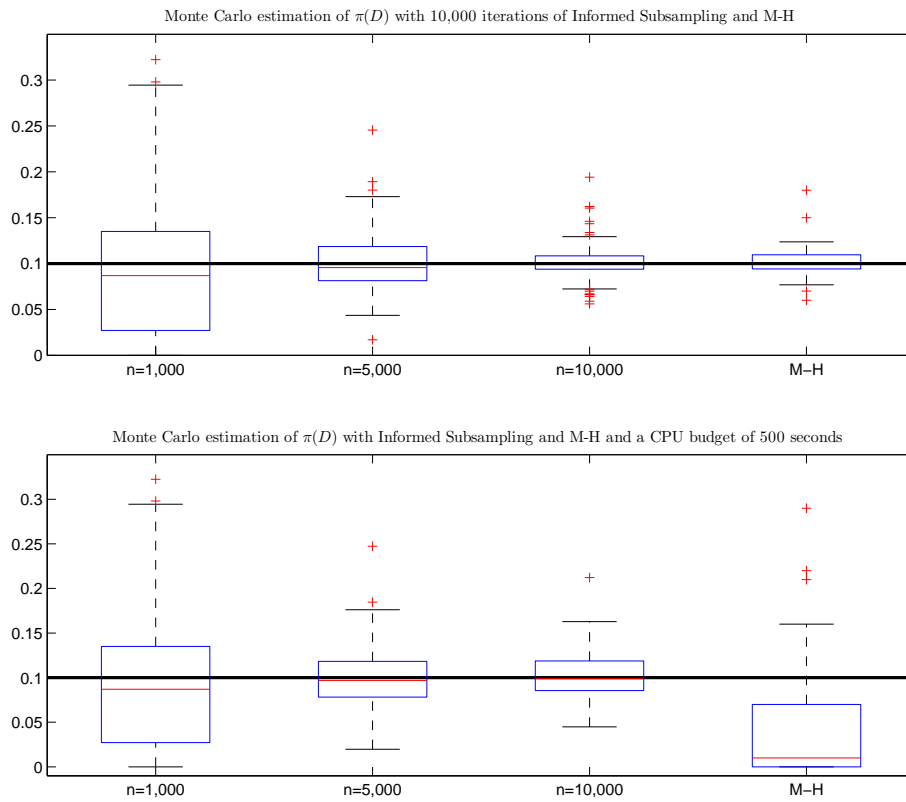


Figure 11: (Example 4: Logistic regression) Estimation of $\pi(D)$ based on ISS-MCMC implemented with $n \in \{1,000; 5,000; 10,000\}$ and Metropolis-Hastings. Top: the experiment is iteration-normalized, *i.e* the chains run for 10,000 iterations. Bottom: the experiment is time-normalized, *i.e* the chains run for 500 seconds. Each chain was replicated 100 times and started from θ^* .

$(\theta_1, \theta_2) \sim_{i.i.d.} p := \mathcal{N}(0, 1/2) \otimes \Gamma(1, 2)$ ($\Gamma(a, b)$ is the Gamma distribution with shape a and rate b) is assigned to θ . Consider an algorithm \mathbf{a} that simulates a Markov chain

$$\theta_{\mathbf{a}} := \left\{ \left(\theta_{1,k}^{(\mathbf{a})}, \theta_{2,k}^{(\mathbf{a})} \right), k \in \mathbb{N} \right\}$$

targeting the posterior distribution of θ given Y , perhaps approximately. We consider the real-time supervised classifier $I_{\mathbf{a}}^*(t)$, driven by $\theta_{\mathbf{a}}$, for the test dataset $Y^* = \{Y_k^*, k \leq N_{test}\}$ ($N_{test} = 10^4$) and defined as:

$$I_{\mathbf{a}}^*(t) = (I_{\mathbf{a},1}^*(t), \dots, I_{\mathbf{a},N_{test}}^*(t)), \quad I_{\mathbf{a},k}^*(t) = \arg \max_{i \in \{1,2\}} f(Y_k^* | \bar{\theta}_{i,\kappa_{\mathbf{a}}(t)}^{(\mathbf{a})}), \quad (36)$$

where $\kappa_{\mathbf{a}}(t) = \sup_{k \in \mathbb{N}} \{\tau_k^{\mathbf{a}} \leq t\}$ and $\bar{\theta}_{i,k}^{(\mathbf{a})} = (1/k) \sum_{\ell=1}^k \theta_{i,\ell}^{(\mathbf{a})}$. We have defined $\tau_k^{\mathbf{a}}$ as the wall clock time to generate k iterations of algorithm \mathbf{a} . We define the live classification error rate as $\epsilon_{\mathbf{a}}(t) = \|I_{\mathbf{a}}^*(t) - I^*\|_1$ where $I_{\mathbf{a}}^* = (I_{\mathbf{a},1}^*, \dots, I_{\mathbf{a},N}^*)$ and I_k^* is the true class of Y_k^* . We compare $\epsilon_{\mathbf{a}}$ for three different algorithms \mathbf{a} : ISS-MCMC, Metropolis-Hastings and Subsampled Likelihoods (Bardenet et al., 2014).

In this simulation example, we have used the true value $\theta^* = (-1, 1/2, 1, 1/2)$ and simulated Y such that it contains the same number of observations from model 1 and model 2 i.e $N/2$. The three algorithms were implemented with the same proposal kernel, namely a single site random walk with adaptive variance that guarantees an acceptance rate between 0.40 and 0.50, see Roberts et al. (2001); Haario et al. (2001). ISS-MCMC was implemented with parameters $n = 1,000$ and $\epsilon = 10^7$. The summary statistics were taken as $S(Y_U) = 0$ if $\sum_{k \in U} \mathbb{1}_{\{I_k=1\}} \neq \sum_{k \in U} \mathbb{1}_{\{I_k=2\}}$ and

$$S(Y_U) = \left[(2/n) \sum_{k=1}^{n/2} Y_k \mathbb{1}_{\{I_k=1\}}, \text{tr}(\text{cov}(Y_k, k \in U, I_k = 1)), \right. \\ \left. (2/n) \sum_{k=1}^{n/2} Y_k \mathbb{1}_{\{I_k=2\}}, \text{tr}(\text{cov}(Y_k, k \in U, I_k = 2)) \right] \quad (37)$$

otherwise. This choice allows to keep the right proportion of data from the two models in any subsample used for the inference. The statistics in (37) are sufficient for each model, taken separately. Subsampled Likelihoods was implemented with the default parameters prescribed in the introduction of Section 4 in Bardenet et al. (2014).

Figure 12 compares the live classification error rate achieved by the three algorithms. We also report the optimal Bayes classifier which achieves $\epsilon_B(t) = 0.0812$ classifying Y_k^* in class 1 if $Y_{k,1}^* < 0$ and in class 2 alternatively. Unsurprisingly, Metropolis-Hastings is penalized because it evaluates the norm of a $N = 10^7$ dimensional vector at each iteration. Subsampled Likelihoods does slightly better than M-H but suffers from the fact that close to stationary regime, the algorithm ends up drawing the quasi-entire dataset with high probability, a fact which was explained in Bardenet et al. (2014).

6.5. Additional details for the handwritten digit inference (Example 1)

In the handwritten digit example (Example 1), we have used batches of $n = 100$ data. Since the initial dataset comprises 2,000 observations per digit, the summary

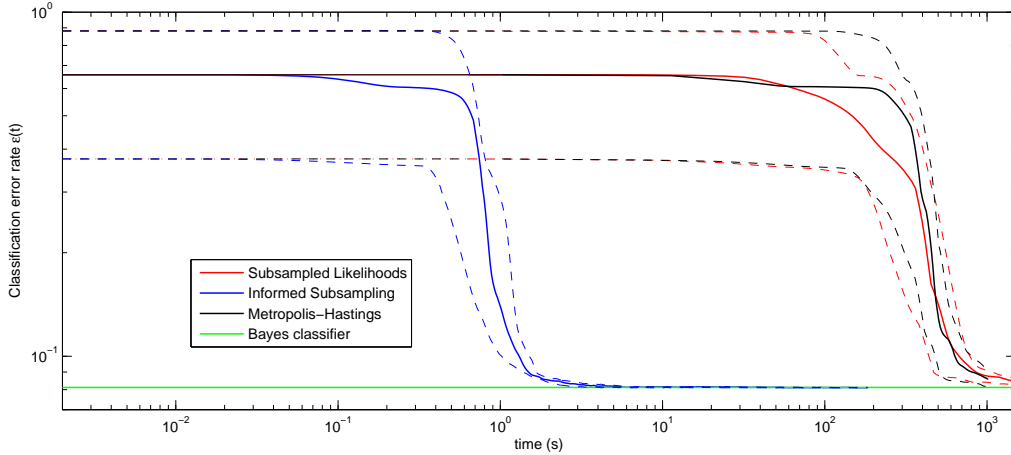


Figure 12: (Example 5: Binary classification) Live classification error rate for three algorithms. This plot was generated by classifying the same test dataset Y^* using the same training dataset Y for the three algorithms. The variability arises from the initial state of the Markov chains. We have used 30 different initial states for the three algorithms and report the median (plain line) and the two quartiles (dashed lines).

statistics were defined in a way that any subsample contains 20 observations from each class. More precisely, we have set for any subsample Y_U , $S(Y_U) = 0$ if for at least one class $i \in \{1, \dots, 5\}$, $(1/n) \sum_{k \in U} \mathbb{1}_{\{J(k)=i\}} \neq 20$ and

$$S(Y_U) = \left\{ \sum_{k \in U} \phi(\theta_{J(k)}) \mathbb{1}_{\{J(k)=i\}} / \sum_{k \in U} \mathbb{1}_{\{J(k)=i\}} \right\}_{i=1}^5$$

otherwise. We set the bandwidth to $\epsilon = 10^5$. The proposal kernels of the Informed Sub-Sampling chain and M-H were defined as the same Random Walk kernel. In particular, at each iteration only a bloc of the template parameter of one of the five classes is updated. The variance parameter of the Random Walk is adapted according to the past trajectory of the chain, so as to maintain an acceptance rate of .25.

Figure 13 reports the empirical marginal distribution of one component for each vector $\theta_1, \dots, \theta_5$ obtained from ISS-MCMC and from M-H. Those distributions are estimated from 50,000 iterations of both algorithms, in stationary regime. This shows that the distribution of those parameters are in line with each other.

7. Conclusion

When the available computational budget is limited, inferring a statistical model based on a tall dataset in the Bayesian paradigm using the Metropolis-Hastings (M-H) algorithm is not computationally efficient. Several variants of the M-H algorithm have been proposed to address this computational issue (Bardenet et al., 2014; Banterle et al., 2015; Korattikara et al., 2014; Maclaurin and Adams, 2015). However, (i) they often

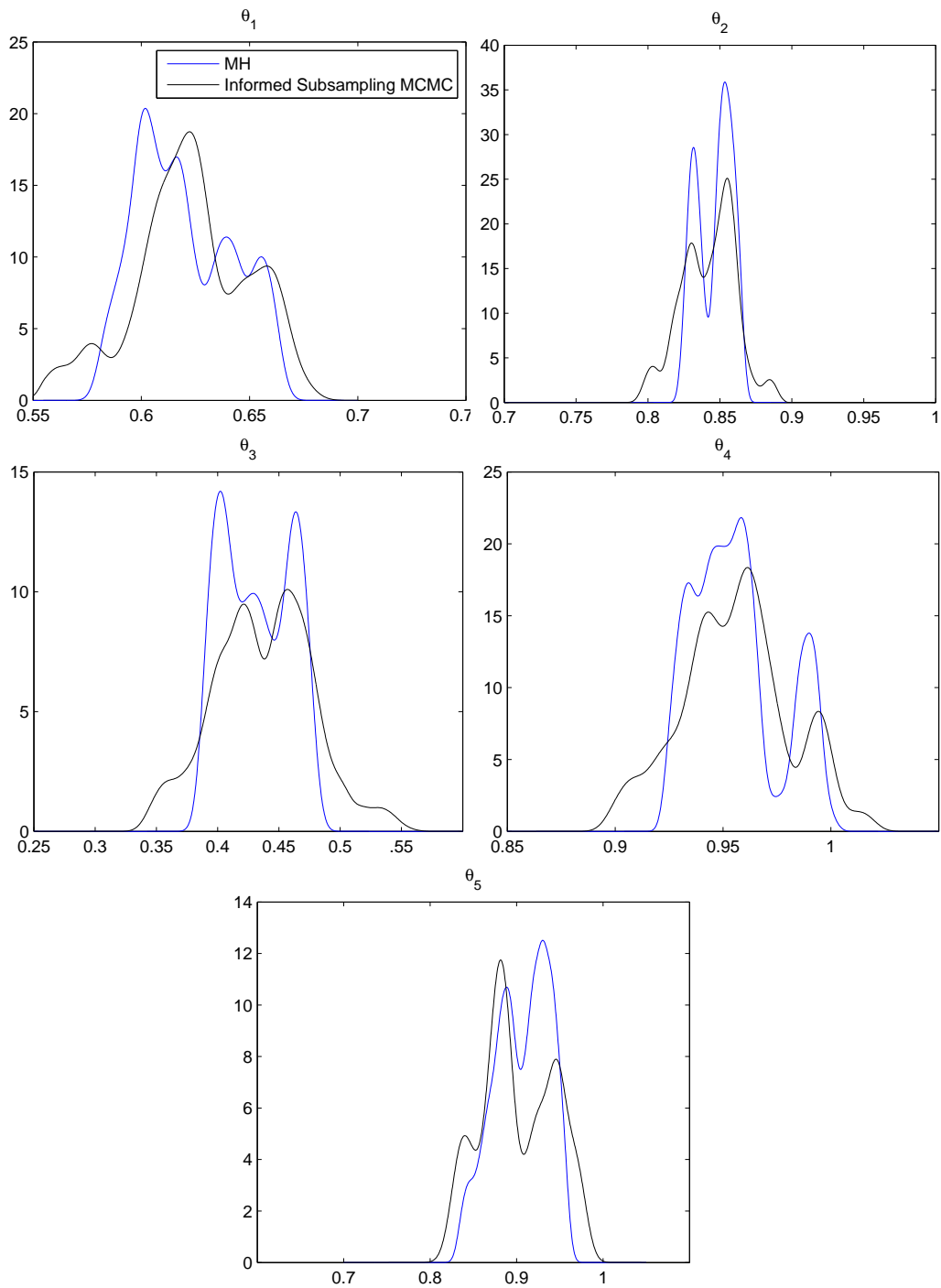


Figure 13: (Example 1: Handwritten digits) Empirical marginal distribution of one component of the vectors $\theta_1, \dots, \theta_5$ using the Metropolis-Hastings chain (blue) and the Informed Sub-Sampling chain (black), estimated from 10,000 transitions at stationary regime.

lose the original simplicity of M–H, (ii) they are only applicable in situations where the data are independent and (iii) the computational cost of one iteration is stochastic which can potentially compromise any computational saving. Informed Sub-Sampling MCMC pushes the approximation one step forward: the computational cost of one iteration is deterministic and is controlled through a user specified parameter, the size of the subsamples.

The aforementioned methods rely on subsampling the whole dataset uniformly at random, at each iteration. Using such subsamples in lieu of the whole dataset in the original M–H chain leads to an algorithm whose approximation of π comes with no guarantee when the subsamples size is fixed. Our main contribution is to show that assigning a distribution to the subsamples that reflects their fidelity to the whole dataset allows to control the L1 approximation error, even when the subsample size is fixed, which is of computational interest.

The algorithm we propose to achieve this task, Informed Sub-Sampling MCMC, offers an alternative to situations where other scalable Metropolis-Hastings variants cannot be implemented (because the model does not satisfy the assumptions *e.g.* independence of the data, existence of a concentration inequality for the model or a cheap lower bound of the likelihood) or are inefficient (because the method ends up using nearly the whole dataset at each iteration). However, in scaling up Metropolis-Hastings there is no free lunch neither. In particular, our method replaces the uniform subsampling approach by a more sophisticated subsampling mechanism involving summary statistics. In this regard, even though our method is in principle widely applicable, it will only be useful in situations where a cheap summary statistics function satisfying Condition 1 is available. In particular, we have shown that our method will give meaningful results when the maximum likelihood estimator is cheap to compute, which somewhat correlates with the optimality of summary statistics in ABC established in [Fearnhead and Prangle \(2012\)](#).

Appendix A. Proofs

Appendix A.1. Proof of Proposition 1

Proof. For notational simplicity and without loss of generality, we take here g as the identity on Θ . Let $n < N$ and U be a subset of $\{1, \dots, N\}$ with cardinal n . Consider the power likelihood:

$$\tilde{f}_n(Y_U | \theta) = f(Y_U | \theta)^{N/n} = \left\{ \prod_{k \in U} f(Y_k | \theta) \right\}^{N/n} = \frac{\exp \left\{ (N/n) \sum_{k \in U} S(Y_k) \right\}^T \theta}{L(\theta)^N},$$

and the corresponding power posterior:

$$\tilde{\pi}_n(\theta | Y_U) = \frac{\exp \left\{ (N/n) \sum_{k \in U} S(Y_k) \right\}^T \theta}{L(\theta)^N} p(\theta) / \tilde{Z}_n(Y_U),$$

where

$$\tilde{Z}_n(Y_U) = \int p(d\theta) \frac{\exp \left\{ (N/n) \sum_{k \in U} S(Y_k) \right\}^T \theta}{L(\theta)^N}.$$

For any θ such that $p(\theta) \neq 0$, write:

$$\log \frac{\pi(\theta | Y_{1:N})}{\tilde{\pi}_n(\theta | Y_U)} = \left\{ \sum_{k=1}^N S(Y_k) - (N/n) \sum_{k \in U} S(Y_k) \right\}^T \theta + \log \frac{\tilde{Z}_n(Y_U)}{Z(Y_{1:N})}. \quad (\text{A.1})$$

and the KL divergence between $\pi(\cdot | Y_{1:N})$ and $\tilde{\pi}(\cdot | Y_U)$, denoted $\text{KL}_n(U)$, simply writes

$$\text{KL}_n(U) = \Delta_n(U)^T \mathbb{E}_\pi(\theta) + \log \frac{\tilde{Z}_n(Y_U)}{Z(Y_{1:N})}, \quad (\text{A.2})$$

where $\Delta_n(U) = \sum_{k=1}^N S(Y_k) - (N/n) \sum_{k \in U} S(Y_k)$. Now, note that

$$\begin{aligned} \tilde{Z}_n(Y_U) &= \int p(d\theta) \frac{\exp \left\{ (N/n) \sum_{k \in U} S(Y_k) \right\}^T \theta}{L(\theta)^N} = \int p(d\theta) \frac{\exp \left\{ \sum_{k=1}^N S(Y_k) - \Delta_n(U) \right\}^T \theta}{L(\theta)^N} \\ &= \int p(d\theta) f(Y_{1:N} | \theta) \exp \left\{ -\Delta_n(U)^T \theta \right\} = Z(Y_{1:N}) \mathbb{E}_\pi \left\{ \exp \left(-\Delta_n(U)^T \theta \right) \right\}. \end{aligned} \quad (\text{A.3})$$

Plugging (B.2) into (B.8) yields:

$$\begin{aligned} \text{KL}_n(U) &= \Delta_n(U)^T \mathbb{E}_\pi(\theta) + \log \mathbb{E}_\pi \left\{ \exp \left(-\Delta_n(U)^T \theta \right) \right\}, \\ &= \log \frac{\mathbb{E}_\pi \left\{ \exp \left(-\Delta_n(U)^T \theta \right) \right\}}{\exp \left(-\Delta_n(U)^T \mathbb{E}_\pi(\theta) \right)} = \log \mathbb{E}_\pi \exp \left[\left\{ \mathbb{E}_\pi(\theta) - \theta \right\}^T \Delta_n(U) \right]. \end{aligned} \quad (\text{A.4})$$

Finally, Cauchy-Schwartz inequality provides the following upper bound for $\text{KL}_n(U)$:

$$\text{KL}_n(U) \leq \log \mathbb{E}_\pi \exp \left\{ \left\| \mathbb{E}_\pi(\theta) - \theta \right\| \left\| \Delta_n(U) \right\| \right\}. \quad (\text{A.5})$$

□

Appendix A.2. Proof of Proposition 2

Proof. Under some weak assumptions, Bernstein-von Mises theorem states that $\pi(\cdot | Y_{1:N})$ is asymptotically (in N) a Gaussian distribution with the maximum likelihood θ^* as mean and $\Gamma_N = I^{-1}(\theta^*)/N$ as covariance matrix, where $I(\theta)$ is the Fisher information matrix at θ . Let us denote by Φ the pdf of $\mathcal{N}(\theta^*, \Gamma_N)$. Under this approximation, $\mathbb{E}_\pi(\theta) = \theta^*$ and from (B.2), we write:

$$\begin{aligned} \exp \text{KL}_n(U) &\approx \int \Phi(d\theta) \exp \left[\left\{ \theta^* - \theta \right\}^T \Delta_n(U) \right] = \int \Phi(\theta^* - \theta) \exp \left\{ \theta^T \Delta_n(U) \right\} d\theta \\ &= \int \frac{1}{(2\pi)^{(d/2)} |\Gamma_N|^{(1/2)}} \exp \left\{ -(1/2) \theta^T \Gamma_N^{-1} \theta + \theta^T \Delta_n(U) \right\} d\theta, \\ &= \frac{1}{(2\pi)^{(d/2)} |\Gamma_N|^{(1/2)}} \int \exp \left[-(1/2) \left\{ \theta^T \Gamma_N^{-1} \theta - 2\theta^T \Gamma_N^{-1} \Gamma_N \Delta_n(U) \right\} \right] d\theta, \\ &= \exp \left\{ (1/2) \Delta_n(U)^T \Gamma_N \Delta_n(U) \right\}, \end{aligned} \quad (\text{A.6})$$

by integration of a multivariate Gaussian density function. Eventually, (B.4) yields the following approximation:

$$\text{KL}_n(U) \approx \widehat{\text{KL}}_n(U) = (1/2) \Delta_n(U)^T \Gamma_N \Delta_n(U). \quad (\text{A.7})$$

□

Appendix A.3. Proof of Proposition 3

Proof. Let $U_n \supset A_n(\theta) := \{U \in U_n, g(\theta)^T \Delta_n(U) \leq 0\}$ and remark that using Cauchy-Schwartz inequality, we have:

$$\mathbb{E} \left\{ \frac{f(Y|\theta)}{f(Y_U|\theta)^{N/n}} \right\} \leq \nu_{n,\epsilon} \{A_n(\theta)\} + \sum_{U \in U_n \setminus A_n(\theta)} \nu_{n,\epsilon}(U) \exp\{\|g(\theta)\| \|\Delta_n(U)\|\}.$$

Now, define $\bar{\Delta}_n(U) := \bar{S}(Y) - \bar{S}(Y_U)$ where \bar{S} is the normalized summary statistics vector, *i.e.* $\bar{\Delta}_n(U) = \Delta_n(U)/N$. Clearly, when $N \rightarrow \infty$, some terms

$$\exp\{\|g(\theta)\| \|\Delta_n(U)\|\} = \exp\{N\|g(\theta)\| \|\bar{\Delta}_n(U)\|\}$$

will have a large contribution to the sum. More precisely, any mismatch between summary statistics of some subsamples $\{Y_U, U \in U_n \setminus A_n(\theta)\}$ with respect to the full dataset will be amplified by the factor N , whereby exponentially inflating the upper bound. However, assigning the distribution $\nu_{n,\epsilon}$ (12) to the subsamples $\{Y_U, U \in U_n\}$, allows to balance out this effect. Indeed, note that

$$\mathbb{E} \left\{ \frac{f(Y|\theta)}{f(Y_U|\theta)^{N/n}} \right\} \leq \nu_{n,\epsilon} \{A_n(\theta)\} + \sum_{U \in U_n \setminus A_n(\theta)} \exp\{-\epsilon \|\Delta_n(U)\|^2 + \|g(\theta)\| \|\Delta_n(U)\|\} / Z(\epsilon),$$

where $Z(\epsilon) = \sum_{U \in U_n} \exp\{-\epsilon \|\Delta_n(U)\|^2\}$ and we have, for a fixed n and when $N \rightarrow \infty$, that

$$\nu_{n,\epsilon}(U) \frac{f(Y|\theta)}{f(Y_U|\theta)} \rightarrow_{\|\Delta_n(U)\| \rightarrow \infty} 0.$$

Since g is bounded, then $\mathbb{E} \{f(Y|\theta)/f(Y_U|\theta)^{N/n}\}$ is bounded too. \square

Appendix A.4. Proof of Proposition 4

This result combines ideas that can be found in proofs (Mitrophanov, 2005, Theorem 3.1) and (Alquier et al., 2016, Corollary 2.3). We first make the two following observations related to the Informed Sub-Sampling Markov chain, in preparation to the application of the noisy MCMC framework to our setting:

- by construction, the conditional distribution of $\{(\tilde{\theta}_i, U_i), 1 \leq i \leq n\}$ given U_0 is

$$\begin{aligned} & p(d\tilde{\theta}_1, U_1, \dots, d\tilde{\theta}_i, U_i | U_0) \\ &= H(U_0, U_1) \int_{\Theta} \mu(d\tilde{\theta}_0) \tilde{K}(\tilde{\theta}_0, d\tilde{\theta}_1 | U_1) \times \dots \times H(U_{i-1}, U_i) \tilde{K}(\tilde{\theta}_{i-1}, d\tilde{\theta}_i | U_i), \end{aligned}$$

where H and \tilde{K} are the kernels related to the transitions $U_i \rightarrow U_{i+1}$ and $\tilde{\theta}_i \rightarrow \tilde{\theta}_{i+1}$ given U_{i+1} , respectively described in steps 3–9 and steps 10–16 of Algorithm 2 and μ is the distribution of $\tilde{\theta}_0$. Integrating out all the $\tilde{\theta}$'s leads to

$$p(U_1, \dots, U_i | U_0) = H(U_0, U_1) \times \dots \times H(U_{i-1}, U_i)$$

and $\{U_i, i \in \mathbb{N}\}$ is clearly a Markov chain with transition kernel H . In addition, note that H is a uniformly ergodic transition kernel that admits $\nu_{n,\epsilon}$ as stationary distribution. We define p_i as the distribution of U_i :

$$p_i(U_i) = \sum_{U_0 \in U_n} p_0(U_0) H^i(U_0, U_i).$$

- the marginal Markov chain $\{\tilde{\theta}_i, i \in \mathbb{N}\}$ produced by our algorithm is non-homogeneous since

$$p(d\tilde{\theta}_i | \tilde{\theta}_{i-1}) = \sum_{U \in \mathcal{U}_n} \tilde{K}(\tilde{\theta}_{i-1}, d\tilde{\theta}_i | U) p_i(U), \quad (\text{A.8})$$

depends on i . We denote by \tilde{K}_i the marginal transition kernel $\tilde{\theta}_{i-1} \rightarrow \tilde{\theta}_i$. However, we observe that as soon as the chain $\{U_i, i \in \mathbb{N}\}$ reaches stationarity, \tilde{K}_i becomes homogeneous as $p_i \rightarrow \nu_{n,\epsilon}$ and is therefore independent of i .

As a consequence, (Alquier et al., 2016, Corollary 2.3) cannot be directly applied to our setting. Indeed, the sequence of auxiliary variables (U_1, U_2, \dots) that appears in the noisy acceptance ratio (19) is not *i.i.d.* but is a Markov chain. We need the following Lemma:

Lemma Appendix A.1. *Let K be the transition kernel of an uniformly ergodic Markov chain that admits π as stationary distribution. Let \tilde{K}_i be the i -th transition kernel of the noisy Markov chain whose auxiliary variable U_i is non independent and non identically distributed. In particular, let p_i be the distribution of the random variable U_i , used at iteration i of the noisy Markov chain. We have:*

$$\lim_{i \rightarrow \infty} \|\pi - \tilde{\pi}_i\|_{\text{TV}} \leq \kappa \sup_{\theta \in \Theta} \sup_{k \in \mathbb{N}} \int \delta_k(\theta, \theta') Q(\theta, d\theta'), \quad (\text{A.9})$$

where $\delta_k : \Theta \times \Theta \rightarrow \mathbb{R}^+$ is a function that satisfies

$$\mathbb{E}_k \{ |\alpha(\theta, \theta') - \tilde{\alpha}(\theta, \theta' | U) | \} \leq \delta_k(\theta, \theta')$$

and the expectation is under p_k .

The proof of this Lemma is postponed to Appendix A.5 and is essentially extending the arguments of Mitrophanov (2005) (Theorem 3.1) and Alquier et al. (2016) (Corollary 2.3) to the case where the noisy Markov chain is non-homogeneous.

By straightforward algebra, we have:

$$\mathbb{E}_k |\alpha(\theta, \theta') - \tilde{\alpha}(\theta, \theta' | U_k)| = \alpha(\theta, \theta') \mathbb{E}_k \left\{ \frac{f(Y | \theta)}{f(Y_U | \theta)^{N/n}} |\phi_U(\theta) - \phi_U(\theta')| \right\} \quad (\text{A.10})$$

where we have defined $\phi_U(\theta) = f(Y_U | \theta) / f(Y | \theta)^{N/n}$. Using Lemma 1, we have that

$$\begin{aligned} \lim_{i \rightarrow \infty} \|\pi - \tilde{\pi}_i\| &\leq \kappa \sup_{\theta \in \Theta} \sup_{k \in \mathbb{N}} \mathbb{E}_k \left\{ \frac{f(Y | \theta)}{f(Y_U | \theta)^{N/n}} \int Q(\theta, d\theta') \alpha(\theta, \theta') |\phi_U(\theta) - \phi_U(\theta')| \right\}, \\ &\leq \kappa \sup_{\theta \in \Theta} \sup_{k \in \mathbb{N}} \mathbb{E}_k \left\{ \frac{f(Y | \theta)}{f(Y_U | \theta)^{N/n}} \right\} \sup_{U \in \mathcal{U}_n} \int Q(\theta, d\theta') \alpha(\theta, \theta') |\phi_U(\theta) - \phi_U(\theta')|. \end{aligned} \quad (\text{A.11})$$

To prove (25), it is sufficient to see that all the expectations in the RHS of (A.11) are equal. Indeed, by construction of the ISS-MCMC algorithm, if $p_0 = \nu_{n,\epsilon}$, then for all $k \in \mathbb{N}$, $p_k = \nu_{n,\epsilon}$.

Appendix A.5. Proof of Lemma 1

In addition of the notations of Section 4, we define the following quantities for a Markov transition kernel regarded as an operator on \mathcal{M} , the space of signed measures on $(\Theta, \mathcal{B}(\Theta))$: $\tau(K) := \sup_{\pi \in \mathcal{M}_{0,1}} \|\pi K\|_{\text{TV}}$ is the ergodicity coefficient of K , $\|K\| := \sup_{\pi \in \mathcal{M}_1} \|\pi K\|$ is the operator norm of K and $\mathcal{M}_1 := \{\pi \in \mathcal{M}, \|\pi\|_{\text{TV}} = 1\}$ and $\mathcal{M}_{0,1} := \{\pi \in \mathcal{M}_1, \pi(\Theta) = 0\}$.

Based on the remarks of Section 4, the assumption that the sequence of auxiliary variable U_1, U_2, \dots is not *i.i.d.* implies that $\{\tilde{\theta}_n, n \in \mathbb{N}\}$ is a non-homogeneous Markov chain with transition kernel $\{\tilde{K}_n, n \in \mathbb{N}\}$. For each $n \in \mathbb{N}$, define π_n as the distribution of θ_n produced by the Metropolis-Hastings algorithm (Alg. 1) with transition kernel K , referred to as the exact kernel hereafter. Our proof is based on the following identity:

$$K^n - \tilde{K}_1 \tilde{K}_2 \cdots \tilde{K}_n = (K - \tilde{K}_1)K^{n-1} + \tilde{K}_1(K - \tilde{K}_2)K^{n-2} + \tilde{K}_1 \tilde{K}_2(K - \tilde{K}_3)K^{n-3} + \cdots + \tilde{K}_1 \cdots \tilde{K}_{n-1}(K - \tilde{K}_n), \quad (\text{A.12})$$

for each $n \in \mathbb{N}$. Equation (A.12) will help translating the proof of Theorem 3.1 in Mitrophanov (2005) to the non-homogeneous setting and in particular, we have for each $n \in \mathbb{N}$:

$$\pi_n - \tilde{\pi}_n = (\pi_0 - \tilde{\pi}_0)K^n + \sum_{i=0}^{n-1} \tilde{\pi}_i(K - \tilde{K}_{i+1})K^{n-i-1}. \quad (\text{A.13})$$

Following the proof of Theorem 3.1 in Mitrophanov (2005), we obtain

$$\begin{aligned} \|\pi_n - \tilde{\pi}_n\|_{\text{TV}} &\leq \|\pi_0 - \tilde{\pi}_0\|_{\text{TV}} \tau(K^n) + \sum_{i=0}^{n-1} \|K - \tilde{K}_{n-i}\| \tau(K^i), \\ &\leq \begin{cases} \|\pi_0 - \tilde{\pi}_0\|_{\text{TV}} + n \sup_{i \leq n} \|K - \tilde{K}_i\| & \text{if } n \leq \lambda \\ \|\pi_0 - \tilde{\pi}_0\|_{\text{TV}} C \rho^n + \sup_{i \leq n} \|K - \tilde{K}_i\| \left\{ \lambda + C \frac{\rho^\lambda - \rho^n}{1 - \rho} \right\} & \text{else} \end{cases} \end{aligned} \quad (\text{A.14})$$

where $\lambda = \lceil \log_\rho(1/C) \rceil$. This yields the following upper bound:

$$\sup_{n \in \mathbb{N}} \|\pi_n - \tilde{\pi}_n\|_{\text{TV}} \leq \|\pi_0 - \tilde{\pi}_0\|_{\text{TV}} + \left(\lambda + C \frac{\rho^\lambda}{1 - \rho} \right) \sup_{n \in \mathbb{N}} \|K - \tilde{K}_n\|. \quad (\text{A.15})$$

Since π is the stationary distribution of $\{\theta_k, k \in \mathbb{N}\}$, (A.14) yields

$$\lim_{n \rightarrow \infty} \|\pi - \pi_n\| \leq \left(\lambda + C \frac{\rho^\lambda}{1 - \rho} \right) \sup_{n \in \mathbb{N}} \|K - \tilde{K}_n\|. \quad (\text{A.16})$$

Using a similar derivation than in the proof of Corollary 2.3 in Alquier et al. (2016), we obtain

$$\|K - \tilde{K}_n\| \leq \sup_{\theta \in \Theta} \int Q(\theta, d\theta') \mathbb{E}_n |\alpha(\theta, \theta') - \tilde{\alpha}(\theta, \theta' | U_n)|,$$

where the expectation is under p_n and which combined with (A.16) leads to

$$\lim_{n \rightarrow \infty} \|\pi_n - \tilde{\pi}_n\|_{\text{TV}} \leq \left(\lambda + C \frac{\rho^\lambda}{1 - \rho} \right) \sup_{\theta \in \Theta} \sup_{n \in \mathbb{N}} \mathbb{E} |\alpha(\theta, \theta') - \tilde{\alpha}(\theta, \theta' | U_n)|$$

where the expectation is under $Q(\theta, \cdot) \otimes p_n$. Any upper bound $\delta_n(\theta, \theta')$ of the expectation on the right hand side yields (A.9).

Appendix A.6. Proof of Proposition 5

Proof. Note that for all $(\theta, \zeta) \in \Theta \times \mathbb{R}^d$, a Taylor expansion of $\pi(\theta)$ and $\phi_U(\theta)$ at $\theta + \Sigma\zeta$ in (26) combined to the triangle inequality leads to:

$$B(U, \theta) \leq \frac{1}{\sqrt{N}} \mathbb{E} \left\{ |(M\zeta)^T \nabla_{\theta} \phi_U(\theta)| \left(1 + \frac{1}{\sqrt{N}} (M\zeta)^T \nabla_{\theta} \log \pi(\theta) \right) \right\} \\ + \frac{1}{2N} \mathbb{E} \{ |(M\zeta)^T \nabla_{\theta}^2 \phi_U(\theta) M\zeta| \} + \mathbb{E} \{ R(\|M\zeta\|/\sqrt{N}) \},$$

where the expectation is under Φ_d and $R(x) = o(x)$ at 0. Applying Cauchy-Schwartz gives:

$$B(U, \theta) \leq \frac{1}{\sqrt{N}} \mathbb{E} \{ \|M\zeta\| \} \|\nabla_{\theta} \phi_U(\theta)\| + \frac{1}{N} \mathbb{E} \{ \|M\zeta\|^2 \} \|\nabla_{\theta} \phi_U(\theta)\| \|\nabla_{\theta} \log \pi(\theta)\| \\ + \frac{1}{2N} \mathbb{E} \{ |\zeta^T M^T \nabla_{\theta}^2 \phi_U(\theta) M\zeta| \} + \mathbb{E} \{ R(\|M\zeta\|/\sqrt{N}) \}.$$

Now, we observe that:

- $\mathbb{E} \{ \|M\zeta\| \} = \mathbb{E} \{ \sum_{i=1}^d (\sum_{j=1}^d M_{i,j} \zeta_j)^2 \}^{1/2} \leq \mathbb{E} \{ \sum_{i=1}^d |\sum_{j=1}^d M_{i,j} \zeta_j| \} \leq \mathbb{E} \{ \sum_{i=1}^d \sum_{j=1}^d |M_{i,j} \zeta_j| \} = \sum_{i=1}^d \sum_{j=1}^d |M_{i,j}| \mathbb{E} \{ |\zeta_i| \} = \sqrt{\frac{2}{\pi}} \|M\|_1$
- $\mathbb{E} \{ \|M\zeta\|^2 \} = \mathbb{E} \{ \sum_{i=1}^d (\sum_{j=1}^d M_{i,j} \zeta_j)^2 \} = \sum_{i=1}^d \mathbb{E} \{ (\sum_{j=1}^d M_{i,j} \zeta_j)^2 \} = \sum_{i=1}^d \text{var}(\sum_{j=1}^d M_{i,j} \zeta_j) = \sum_{i=1}^d \sum_{j=1}^d M_{i,j}^2 \text{var}(\zeta_j) = \|M\|_2^2$
- considering the quadratic form associate to the operator $T(U, \theta) = M^T \nabla_{\theta}^2 \phi_U(\theta) M$, noting that $T(U, \theta)$ is symmetric its eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ are real and we have

$$\zeta^T T(U, \theta) \zeta \leq \lambda_1 \|\zeta\|^2$$

so that:

$$\mathbb{E} \{ |(M\zeta)^T \nabla_{\theta}^2 \phi_U(\theta) M\zeta| \} \leq d \sup_i |\lambda_i| \leq d \|M^T \nabla_{\theta}^2 \phi_U(\theta) M\|$$

where for any square matrix A , we have defined $\|A\| = \sup_{x \in \mathbb{R}^d, \|x\|=1} \|Ax\|$ as the operator norm. □

Appendix B. Proof of Proposition 6

In this section, we are assuming that there is an infinite stream of observations (Y_1, Y_2, \dots) and a parameter $\theta_0 \in \Theta$ such that $Y_i \sim f(\cdot | \theta_0)$. Let $\rho > 1$ be a constant defined as the ratio N/n *i.e* the size of the full dataset over the size of the subsamples of interest. The full dataset is thus $Y_{1:\rho n}$. We define the set

$$\mathbf{U}_n^{\rho} = \{U \subset \{1, \dots, \rho n\}, |U| = n\}$$

such that Y_U ($U \in \mathbf{U}_n^{\rho}$) is the set of subsamples of interest. We study the asymptotics when $n \rightarrow \infty$ *i.e* we let the whole dataset and the size of subsamples of interest grow at the same rate.

Proposition 6. Let $\theta_{\rho n}^*$ be the MLE of $Y_1, \dots, Y_{\rho n}$ and θ_U^* be the MLE of the subsample Y_U ($U \in \mathcal{U}_n^\rho$). Assume that there exists a compact set $\kappa_n \subset \Theta$ such that $(\theta_{\rho n}^*, \theta_0) \in \kappa_n^2$ and for all U , there exists a compact set $\kappa_U \subset \Theta$ such that $(\theta_U^*, \theta_0) \in \kappa_U^2$. Then, there exists a constants β , a metric $\|\cdot\|_{\theta_0}$ on Θ and a non-decreasing subsequence $\{\sigma_n\}_{n \in \mathbb{N}}$, ($\sigma_n \in \mathbb{N}$) such that for all $U \in \mathcal{U}_{\sigma_n}^\rho$, we have for p -almost all $\theta \in \kappa_n \cap \kappa_U$

$$\log f(Y_{1:\rho\sigma_n} | \theta) - \rho \log f(Y_U | \theta) \leq H_n(Y, \theta) + \beta + \frac{\rho\sigma_n}{2} \|\theta_U^* - \theta^*\|_{\theta_0}, \quad (\text{B.1})$$

where

$$\text{plim}_{n \rightarrow \infty} H_n(Y, \theta) \stackrel{\mathbb{P}_{\theta_0}}{=} 0.$$

Proof. Fix $n \in \mathbb{N}$. Consider the case where the prior distribution p is uniform on κ_n . In this case, the posterior is

$$\pi_n(\theta | Y_{1:\rho n}) = f(Y_{1:\rho n} | \theta) \mathbb{1}_{\kappa_n}(\theta) / Z_{\rho n}, \quad Z_{\rho n} = \int_{\kappa_n} f(Y_{1:\rho n} | \theta) d\theta$$

and from corollary 2, we know that there exists a subsequence $\tau_n \subset \mathbb{N}$ such that for p -almost all $\theta \in \kappa_n$

$$\left| \log \frac{f(Y_{1:\rho\tau_n} | \theta)}{Z_{\rho\tau_n}} - \log \Phi_{\rho\tau_n}(\theta) \right| \stackrel{\mathbb{P}_{\theta_0}}{\rightarrow} 0, \quad (\text{B.2})$$

where $\theta \mapsto \Phi_{\rho\tau_n}(\theta)$ is the pdf of $\mathcal{N}(\theta_{\rho\tau_n}^*, I(\theta_0)^{-1}/\rho\tau_n)$. Similarly, there exists another subsequence $\gamma_n \subset \mathbb{N}$ such that for all $U \in \mathcal{U}_{\gamma_n}^\rho$ and for p -almost all $\theta \in \kappa_U$

$$\left| \rho \log \frac{f(Y_U | \theta)}{Z_{\gamma_n}(U)} - \rho \log \Phi_U(\theta) \right| \stackrel{\mathbb{P}_{\theta_0}}{\rightarrow} 0, \quad Z_{\gamma_n}(U) = \int_{\kappa_U} f(Y_U | \theta) d\theta \quad (\text{B.3})$$

where $\theta \mapsto \Phi_U(\theta)$ is the pdf of $\mathcal{N}(\theta_U^*, I(\theta_0)^{-1}/|U|)$. Let $\{\sigma_n\}_{n \in \mathbb{N}}$ be the sequence defined as $\sigma_n = \max\{\tau_n, \gamma_n\}$. We know from (B.2) and (B.3) that for all $\varepsilon > 0$ and all $\eta > 0$, there exists $n_1 \in \mathbb{N}$ such that for all $U \in \mathcal{U}_{\sigma_n}^\rho$ and for all $n \geq n_1$

$$\mathbb{P}_{\theta_0} \left\{ \left| \log \frac{f(Y_{1:\rho\sigma_n} | \theta)}{Z_{\rho\sigma_n}} - \log \Phi_{\rho\sigma_n}(\theta) \right| + \left| \rho \log \frac{f(Y_U | \theta)}{Z_{\sigma_n}(U)} - \rho \log \Phi_U(\theta) \right| \geq \varepsilon \right\} \leq \eta. \quad (\text{B.4})$$

Now, by straightforward algebra, we have for any $U \in \mathcal{U}_{\sigma_n}^\rho$

$$\begin{aligned} \log f(Y_{1:\rho\sigma_n} | \theta) - \rho \log f(Y_U | \theta) &= \log \frac{f(Y_{1:\rho\sigma_n} | \theta)}{Z_{\rho\sigma_n}} - \log \Phi_{\rho\sigma_n}(\theta) - \rho \log \frac{f(Y_U | \theta)}{Z_{\sigma_n}(U)} \\ &\quad + \rho \log \Phi_U(\theta) + \log \frac{Z_{\rho\sigma_n}}{Z_{\sigma_n}(U)^\rho} + \log \Phi_{\rho\sigma_n}(\theta) - \rho \log \Phi_U(\theta) \\ &\leq \left| \log \frac{f(Y_{1:\rho\sigma_n} | \theta)}{Z_{\rho\sigma_n}} - \log \Phi_{\rho\sigma_n}(\theta) - \rho \log \frac{f(Y_U | \theta)}{Z_{\sigma_n}(U)} + \rho \log \Phi_U(\theta) \right| \\ &\quad + \log \frac{Z_{\rho\sigma_n}}{Z_{\sigma_n}(U)^\rho} + (\rho - 1) \log(2\pi)^{d/2} + \frac{\rho\sigma_n}{2} \left| \|\theta - \theta_U^*\|_{\theta_0} - \|\theta - \theta^*\|_{\theta_0} \right| \\ &\leq \left| \log \frac{f(Y_{1:\rho\sigma_n} | \theta)}{Z_{\rho\sigma_n}} - \log \Phi_{\rho\sigma_n}(\theta) \right| + \left| \rho \log \frac{f(Y_U | \theta)}{Z_{\sigma_n}(U)} - \rho \log \Phi_U(\theta) \right| \\ &\quad + \log \frac{Z_{\rho\sigma_n}}{Z_{\sigma_n}(U)^\rho} + (\rho - 1) \log(2\pi)^{d/2} + \frac{\rho\sigma_n}{2} \|\theta_U^* - \theta^*\|_{\theta_0}, \quad (\text{B.5}) \end{aligned}$$

where we have used Lemma 2 for the first inequality and the triangle inequalities for the second. Combining (B.5) with (B.4) yields (B.1). \square

Lemma 1. *Consider a posterior distribution π_n given n data $Y_{1:n}$ where p is the prior distribution and its Bernstein-von Mises approximation is $\Phi_n = \mathcal{N}(\theta^*(Y_{1:n}), I(\theta_0)^{-1}/n)$. There exists a subsequence $\{\tau_n\}_n \subset \mathbb{N}$ such that*

$$\text{plim}_{n \rightarrow \infty} |\pi_{\tau_n}(\theta) - \Phi_{\tau_n}(\theta)| \stackrel{\mathbb{P}_{\theta_0}}{=} 0, \quad \text{for } p\text{-almost all } \theta. \quad (\text{B.6})$$

Proof. This follows for the fact that convergence in L_1 implies pointwise convergence almost everywhere of a subsequence, i.e there exists a subsequence $\{\tau_n\}_{n \in \mathbb{N}} \subset \mathbb{N}$ such that

$$\|\pi_n - \Phi_n\|_1 \rightarrow 0 \Rightarrow |\pi_{\tau_n}(\theta) - \Phi_{\tau_n}(\theta)| \rightarrow 0 \quad p\text{-a.e.} \quad (\text{B.7})$$

Eq. B.6 follows from combining the Bernstein-von Mises theorem and Eq. (B.7):

$$\text{plim}_{n \rightarrow \infty} \|\pi_n - \Phi_n(\theta^*, I(\theta_0)^{-1}/n)\|_1 \stackrel{\mathbb{P}_{\theta_0}}{=} 0 \Rightarrow \text{plim}_{n \rightarrow \infty} |\pi_{\tau_n}(\theta) - \Phi_{\tau_n}(\theta)| \stackrel{\mathbb{P}_{\theta_0}}{=} 0 \quad p\text{-a.e.}$$

\square

Corollary 2. *There exists a subsequence $\{\tau_n\}_{n \in \mathbb{N}} \subset \mathbb{N}$ such that*

$$\text{plim}_{n \rightarrow \infty} |\log \pi_{\tau_n}(\theta) - \log \Phi_{\tau_n}(\theta)| \stackrel{\mathbb{P}_{\theta_0}}{=} 0, \quad \text{for } p\text{-almost all } \theta. \quad (\text{B.8})$$

Proof. Follows from Lemma 1, by continuity of the logarithm. \square

Lemma 2. *For any $U \in \mathcal{U}_n$, let $\theta \mapsto \Phi_U(\theta)$ be the pdf of $\mathcal{N}(\theta_U^*, I(\theta_0)^{-1}/n)$ and $\Phi_{\rho n}$ be the pdf of $\mathcal{N}(\theta_{\rho n}^*, I(\theta_0)^{-1}/\rho n)$ be the Bernstein-von Mises approximations of respectively $\pi(\cdot | Y_U)$ and $\pi(\cdot | Y_{1:\rho n})$ where $U \subset \mathcal{U}_n(Y_{1:\rho n})$. Then we have for all $\theta \in \Theta$*

$$\log \Phi_{\rho n}(\theta) - \rho \log \Phi_U(\theta) \leq (\rho - 1) \log(2\pi)^{d/2} + \frac{\rho n}{2} \{ \|\theta - \theta_U^*\|_{\theta_0} - \|\theta - \theta^*\|_{\theta_0} \},$$

where for any d -squared symmetric matrix M , we have defined by $\|\cdot\|_M$ the norm associated to the scalar product $\langle u, v \rangle_M = u^T M v$.

Proof. This follows from straightforward algebra and noting that

$$\log \rho n |I(\theta_0)| - \rho \log n |I(\theta_0)| \leq 0.$$

\square

Acknowledgements

The Insight Centre for Data Analytics is supported by Science Foundation Ireland under Grant Number SFI/12/RC/2289. Nial Friel's research was also supported by an Science Foundation Ireland grant: 12/IP/1424. Pierre Alquier's research was funded by Labex ECODEC (ANR - 11-LABEX-0047) and by the research programme New Challenges for New Data from LCL and GENES, hosted by the Fondation du Risque.

References

References

- Allasonnière, S., Amit, Y., Trouvé, A., 2007. Towards a coherent statistical framework for dense deformable template estimation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 69 (1), 3–29.
- Alquier, P., Friel, N., Everitt, R., Boland, A., 2016. Noisy Monte Carlo: Convergence of Markov chains with approximate transition kernels. *Statistics and Computing* 26 (1-2), 29–47.
- Andrieu, C., Roberts, G. O., 2009. The pseudo-marginal approach for efficient Monte Carlo computations. *The Annals of Statistics*, 697–725.
- Andrieu, C., Vihola, M., 2015. Convergence properties of pseudo-marginal markov chain monte carlo algorithms. *The Annals of Applied Probability* 25 (2), 1030–1077.
- Banterle, M., Grazian, C., Lee, A., Robert, C. P., 2015. Accelerating Metropolis-Hastings algorithms by delayed acceptance. arXiv preprint arXiv:1503.00996.
- Bardenet, R., Doucet, A., Holmes, C., 2014. Towards scaling up Markov chain Monte Carlo: an adaptive subsampling approach. In: *ICML*. pp. 405–413.
- Bardenet, R., Doucet, A., Holmes, C., 2015. On Markov chain Monte Carlo methods for tall data. arXiv preprint arXiv:1505.02827.
- Bierkens, J., Fearnhead, P., Roberts, G., 2016. The Zig-Zag process and super-efficient sampling for Bayesian analysis of big data. arXiv preprint arXiv:1607.03188.
- Chib, S., Greenberg, E., 1995. Understanding the Metropolis-Hastings algorithm. *The American statistician* 49 (4), 327–335.
- Csilléry, K., Blum, M. G., Gaggiotti, O. E., François, O., 2010. Approximate Bayesian computation (ABC) in practice. *Trends in ecology & evolution* 25 (7), 410–418.
- Dalalyan, A. S., 2017. Further and stronger analogy between sampling and optimization: Langevin Monte Carlo and gradient descent. arXiv preprint arXiv:1704.04752.
- Fearnhead, P., Bierkens, J., Pollock, M., Roberts, G. O., 2016. Piecewise deterministic Markov processes for continuous-time Monte Carlo. arXiv preprint arXiv:1611.07873.
- Fearnhead, P., Prangle, D., 2012. Constructing summary statistics for approximate Bayesian computation: semi-automatic approximate Bayesian computation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 74 (3), 419–474.
- Geyer, C. J., Thompson, E. A., 1995. Annealing Markov chain Monte Carlo with applications to ancestral inference. *Journal of the American Statistical Association* 90 (431), 909–920.
- Haario, H., Saksman, E., Tamminen, J., 2001. An adaptive Metropolis algorithm. *Bernoulli*, 223–242.
- Huggins, J., Zou, J., 2016. Quantifying the accuracy of approximate diffusions and Markov chains. In: *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, PLMR*. Vol. 54. pp. 382–391.
- Jacob, P. E., Thiery, A. H., et al., 2015. On nonnegative unbiased estimators. *The Annals of Statistics* 43 (2), 769–784.
- Korattikara, A., Chen, Y., Welling, M., 2014. Austerity in MCMC land: Cutting the Metropolis-Hastings budget. In: *Proceedings of the 31st International Conference on Machine Learning*.
- Le Cam, L., 1953. On some asymptotic properties of maximum likelihood estimates and related Bayes' estimates. *Univ. Calif. Publ. in Statist.* 1, 277–330.
- Le Cam, L., 1986. *Asymptotic methods in statistical decision theory*. Springer Science & Business Media.
- Maclaurin, D., Adams, R. P., 2015. Firefly Monte Carlo: Exact MCMC with subsets of data. In: *Twenty-Fourth International Joint Conference on Artificial Intelligence*.
- Marin, J.-M., Pudlo, P., Robert, C. P., Ryder, R. J., 2012. Approximate Bayesian computational methods. *Statistics and Computing* 22 (6), 1167–1180.
- Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., Teller, E., 1953. Equation of state calculations by fast computing machines. *The journal of chemical physics* 21 (6), 1087–1092.
- Mitrophanov, A. Y., 2005. Sensitivity and convergence of uniformly ergodic Markov chains. *Journal of Applied Probability*, 1003–1014.
- Nunes, M. A., Balding, D. J., 2010. On optimal selection of summary statistics for approximate Bayesian computation. *Statistical applications in genetics and molecular biology* 9 (1).
- Pollock, M., Fearnhead, P., Johansen, A. M., Roberts, G. O., 2016. The scalable Langevin exact algorithm: Bayesian inference for big data. arXiv preprint arXiv:1609.03436.

- Pritchard, J. K., Seielstad, M. T., Perez-Lezaun, A., Feldman, M. W., 1999. Population growth of human Y chromosomes: a study of Y chromosome microsatellites. *Molecular biology and evolution* 16 (12), 1791–1798.
- Quiroz, M., Villani, M., Kohn, R., 2015. Speeding up MCMC by efficient data subsampling. *Riksbank Research Paper Series* (121).
- Quiroz, M., Villani, M., Kohn, R., 2016. Exact subsampling MCMC. *arXiv preprint arXiv:1603.08232*.
- Roberts, G. O., Rosenthal, J. S., et al., 2001. Optimal scaling for various Metropolis-Hastings algorithms. *Statistical science* 16 (4), 351–367.
- Rudolf, D., Schweizer, N., 2015. Perturbation theory for Markov chains via Wasserstein distance. *arXiv preprint arXiv:1503.04123*.
- Van der Vaart, A. W., 2000. *Asymptotic statistics*. Vol. 3. Cambridge university press.
- Wilkinson, R. D., 2013. Approximate Bayesian computation (ABC) gives exact results under the assumption of model error. *Statistical applications in genetics and molecular biology* 12 (2), 129–141.