

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/356601649>

Markov Chain Score Ascent: A Unifying Framework of Variational Inference with Markovian Gradients

Preprint · June 2022

CITATIONS

0

READS

293

5 authors, including:



Kyurae Kim

University of Liverpool

11 PUBLICATIONS 34 CITATIONS

SEE PROFILE

MARKOV CHAIN SCORE ASCENT: A Unifying Framework of Variational Inference with Markovian Gradients

Kyurae Kim
Sogang University
msca8h@sogang.ac.kr

Jisu Oh
Sogang University
jisuhoh@sogang.ac.kr

Jacob R. Gardner
University of Pennsylvania
jacobrg@seas.upenn.edu

Adjji Boussou Dieng
Princeton University
adjji@princeton.edu

Hongseok Kim*
Sogang University
hongseok@sogang.ac.kr

Abstract

Minimizing the inclusive Kullback-Leibler (KL) divergence with stochastic gradient descent (SGD) is challenging since its gradient is defined as an integral over the posterior. Recently, multiple methods have been proposed to run SGD with *biased* gradient estimates obtained from a Markov chain. This paper provides the first non-asymptotic convergence analysis of these methods by establishing their mixing rate and gradient variance. To do this, we demonstrate that these methods—which we collectively refer to as Markov chain score ascent (MCSA) methods—can be cast as special cases of the Markov chain gradient descent framework. Furthermore, by leveraging this new understanding, we develop a novel MCSA scheme, *parallel* MCSA (pMCSA), that achieves a tighter bound on the gradient variance. We demonstrate that this improved theoretical result translates to superior empirical performance.

1 Introduction

Bayesian inference aims to analyze the posterior distribution of an unknown latent variable \mathbf{z} from which data \mathbf{x} is observed. By assuming a model $p(\mathbf{x} | \mathbf{z})$, the posterior $\pi(\mathbf{z})$ is given by Bayes' rule such that $\pi(\mathbf{z}) \propto p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})$ where $p(\mathbf{z})$ represents our prior belief on \mathbf{z} . Instead of working directly with π , variational inference (VI, Blei *et al.* 2017) seeks a *variational approximation* $q(\mathbf{z}; \lambda)$ that is the most similar to π according to a discrepancy measure $d(\pi, q(\cdot; \lambda))$.

The apparent importance of choosing the right discrepancy measure has led to a quest spanning a decade (Dieng *et al.*, 2017; Geffner & Domke, 2021b; Hernandez-Lobato *et al.*, 2016; Li & Turner, 2016; Regli & Silva, 2018; Ruiz & Titsias, 2019; Salimans *et al.*, 2015; Wan *et al.*, 2020; Wang *et al.*, 2018; Zhang *et al.*, 2021). So far, the exclusive (or reverse, backward) Kullback-Leibler (KL) divergence $d_{\text{KL}}(q(\cdot; \lambda) \| \pi)$ has seen “exclusive” use, partly because it is defined as an integral over $q(\mathbf{z}; \lambda)$, which can be approximated efficiently. In contrast, the *inclusive* (or forward) KL is defined as an integral over π as

$$d_{\text{KL}}(\pi \| q(\cdot; \lambda)) = \int \pi(\mathbf{z}) \log \frac{\pi(\mathbf{z})}{q(\mathbf{z}; \lambda)} d\mathbf{z} = \mathbb{E}_{\mathbf{z} \sim \pi(\cdot)} \left[\log \frac{\pi(\mathbf{z})}{q(\mathbf{z}; \lambda)} \right].$$

Since our goal is to approximate π with $q(\cdot; \lambda)$ but the inclusive KL involves an integral over π , we end up facing a chicken-and-egg problem. Despite this challenge, the inclusive KL has con-

*Corresponding author.

sistently drawn attention due to its statistical properties, such as better uncertainty estimates due to its mass covering property (MacKay, 2001; Minka, 2005; Trippe & Turner, 2017).

Recently, Naesseth *et al.* (2020); Ou & Song (2020) have respectively proposed Markovian score climbing (MSC) and joint stochastic approximation (JSA). These methods minimize the inclusive KL using stochastic gradient descent (SGD, Robbins & Monro 1951), where the gradients are estimated using a Markov chain. The Markov chain kernel $K_{\lambda_t}(\mathbf{z}_t, \cdot)$ is π -invariant (Robert & Casella, 2004) and is chosen such that it directly takes advantage of the current variational approximation $q(\cdot; \lambda_t)$. Thus, the quality of the gradients improves over time as the KL divergence decreases. Still, the gradients are non-asymptotically biased and Markovian across adjacent iterations, which sharply contrasts MSC and JSA from classical black-box VI (Kucukelbir *et al.*, 2017; Ranganath *et al.*, 2014), where the gradients are unbiased and independent. While Naesseth *et al.* (2020) has shown the convergence of MSC through the work of Gu & Kong (1998), this result is only asymptotic and does not provide practical insight into the performance of MSC.

In this paper, we address these theoretical gaps by casting MSC and JSA into a general framework we call Markov chain score ascent (MCSA), which we show is a special case of Markov chain gradient descent (MCGD, Duchi *et al.* 2012). This enables the application of the non-asymptotic convergence results of MCGD (Debavelaere *et al.*, 2021; Doan *et al.*, 2020a,b; Duchi *et al.*, 2012; Karimi *et al.*, 2019; Sun *et al.*, 2018; Xiong *et al.*, 2021). For MCGD methods, the fundamental properties affecting the convergence rate are the ergodic convergence rate (ρ) of the MCMC kernel and the gradient variance (G). We analyze ρ and G of MSC and JSA, enabling their practical comparison given a fixed computational budget (N). Furthermore, based on the recent insight that the mixing rate does not affect the convergence rate of MCGD Doan *et al.* (2020a,b), we propose a novel scheme, parallel MCSA (pMCSA), which achieves lower variance by trading off the mixing rate. We verify our theoretical analysis through numerical simulations and compare MSC, JSA, and pMCSA on general Bayesian inference problems. Our experiments show that our proposed method outperforms previous MCSA approaches.

Contribution Summary

1. **Section 4:** We provide the first non-asymptotic theoretical analysis of two recently proposed inclusive KL minimization methods, MSC (**Theorem 1**) and JSA (**Theorem 2**).
2. **Section 3:** To do this, we show that both methods can be viewed as what we call “Markov chain score ascent” (MCSA) methods, which are a special case of MCGD (**Proposition 1**).
3. **Section 5:** In light of this, we develop a novel MCSA method which we call parallel MCSA (pMCSA) that achieves lower gradient variance (**Theorem 3**).
4. **Section 6:** We demonstrate that the improved theoretical performance of pMCSA translates to superior empirical performance across a variety of Bayesian inference tasks.

2 Background

2.1 Inclusive Kullback-Leibler Minimization with Stochastic Gradients

VI with SGD The goal of VI is to find the optimal variational parameter λ identifying $q(\cdot; \lambda) \in \mathcal{Q}$ that minimizes some discrepancy measure $D(\pi, q(\cdot; \lambda))$. A typical way to perform VI is to use stochastic gradient descent (SGD, Robbins & Monro 1951), provided that *unbiased* gradient estimates of the optimization target $\mathbf{g}(\lambda)$ are available such that we can repeat the update

$$\lambda_t = \lambda_{t-1} - \gamma_t \mathbf{g}(\lambda_{t-1})$$

where $\gamma_1, \dots, \gamma_T$ is a step-size schedule.

Inclusive KL Minimization with SGD For inclusive KL minimization, \mathbf{g} should be set as

$$\mathbf{g}(\lambda) = \nabla_\lambda d_{\text{KL}}(\pi \parallel q(\cdot; \lambda)) = \nabla_\lambda \mathbb{H}[\pi, q(\cdot; \lambda)] = -\mathbb{E}_{\mathbf{z} \sim \pi(\cdot)} [\mathbf{s}(\lambda; \mathbf{z})]$$

where $\mathbb{H}[\pi, q(\cdot; \lambda)]$ is the cross-entropy between π and $q(\cdot; \lambda)$, which shows the connection with cross-entropy methods (de Boer *et al.*, 2005), and $\mathbf{s}(\lambda; \mathbf{z}) = \nabla_\lambda \log q(\mathbf{z}; \lambda)$ is known as the *score gradient*. Since inclusive KL minimization with SGD is equivalent to ascending towards the direction of the score, Naesseth *et al.* (2020) coined the term *score climbing*. To better conform with the optimization literature, we instead call this approach *score ascent* as in gradient ascent.

2.2 Markov Chain Gradient Descent

Overview of MCGD Markov chain gradient descent (MCGD, Duchi *et al.* 2012; Sun *et al.* 2018) is a family of algorithms that optimize a function f defined as $f(\lambda) = \int f(\lambda, \eta) \Pi(d\eta)$ where $\Pi(d\eta)$ is a probability measure. MCGD repeats the steps

$$\lambda_{t+1} = \lambda_t - \gamma_t \mathbf{g}(\lambda_t, \eta_t), \quad \eta_t \sim P_{\lambda_{t-1}}(\eta_{t-1}, \cdot) \quad (1)$$

where $P_{\lambda_{t-1}}$ is a Π -invariant Markov chain kernel that may depend on λ_{t-1} . The noise of the gradient is Markovian and non-asymptotically biased, departing from vanilla SGD. Non-asymptotic convergence of this general algorithm has recently started to gather attention as by Debavelaere *et al.* (2021); Doan *et al.* (2020a,b); Duchi *et al.* (2012); Karimi *et al.* (2019); Sun *et al.* (2018).

Applications of MCGD MCGD encompasses an extensive range of problems, including distributed optimization (Ram *et al.*, 2009), reinforcement learning (Doan *et al.*, 2020a; Tadić & Doucet, 2017; Xiong *et al.*, 2021), and expectation-minimization (Karimi *et al.*, 2019), to name a few. This paper extends this list with inclusive KL VI through the MCSA framework.

3 Markov Chain Score Ascent

First, we develop Markov chain score ascent (MCSA), a framework for inclusive KL minimization with MCGD. This framework will establish the connection between MSC/JSA and MCGD.

3.1 Markov Chain Score Ascent as a Special Case of Markov Chain Gradient Descent

As shown in Equation (1), the basic ingredients of MCGD are the target function $f(\lambda, \eta)$, the gradient estimator $\mathbf{g}(\lambda, \eta)$, and the Markov chain kernel $P_\lambda(\eta, \cdot)$. Obtaining MCSA from MCGD boils down to designing \mathbf{g} and P_λ such that $f(\lambda) = d_{\text{KL}}(\pi \parallel q(\cdot; \lambda))$. The following proposition provides sufficient conditions on \mathbf{g} and P_λ to achieve this goal.

Proposition 1. Let $\eta = (\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(N)})$ and a Markov chain kernel $P_\lambda(\eta, \cdot)$ be Π -invariant where Π is defined as

$$\Pi(\eta) = \pi(\mathbf{z}^{(1)}) \pi(\mathbf{z}^{(2)}) \times \dots \times \pi(\mathbf{z}^{(N)}).$$

Then, by defining the target function f and the gradient estimator \mathbf{g} to be

$$f(\lambda, \eta) = -\frac{1}{N} \sum_{n=1}^N \log q(\mathbf{z}^{(n)}; \lambda) - \mathbb{H}[\pi], \quad \text{and} \quad \mathbf{g}(\lambda, \eta) = -\frac{1}{N} \sum_{n=1}^N \mathbf{s}(\mathbf{z}^{(n)}; \lambda)$$

where $\mathbb{H}[\pi]$ is the entropy of π , MCGD results in inclusive KL minimization as

$$\mathbb{E}_\Pi[f(\lambda, \eta)] = d_{\text{KL}}(\pi \parallel q(\cdot; \lambda)), \quad \text{and} \quad \mathbb{E}_\Pi[\mathbf{g}(\lambda, \eta)] = \nabla_\lambda d_{\text{KL}}(\pi \parallel q(\cdot; \lambda)).$$

Proof. See page 17 of the *supplementary material*.

This simple connection between MCGD and VI paves the way toward the non-asymptotic analysis of JSA and MSC. Note that N here can be regarded as the computational budget of each MCGD iteration since the cost of (i) generating the Markov chain samples $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(N)}$ and (ii) computing the gradient \mathbf{g} will linearly increase with N .

In addition, the MCGD framework often assumes P to be geometrically ergodic. An exception is the analysis of Debavelaere *et al.* (2021) where they work with polynomially ergodic kernels.

Assumption 1. (Markov chain kernel) The Markov chain kernel P is geometrically ergodic as

$$d_{\text{TV}}(P_\lambda^n(\eta, \cdot), \Pi) \leq C \rho^n$$

for some positive constant C .

3.2 Non-Asymptotic Convergence of Markov Chain Score Ascent

Non-Asymptotic Convergence Through Proposition 1 and Assumption 1 and some technical assumptions on the objective function, we can apply the existing convergence results of MCGD to MCSA. Table 1 provides a list of relevant results. Apart from properties of the objective function (such as Lipschitz smoothness), the convergence rates are stated in terms of the gradient bound G , kernel mixing rate ρ , and the number of MCGD iteration T . We focus on G and ρ as they are closely related to the design choices of different MCSA algorithms.

Table 1: Convergence Rates of MCGD Algorithms

Algorithm	Stepsize Rule	Gradient Assumption	Rate	Reference
Mirror Descent ¹	$\gamma_t = \gamma/\sqrt{t}$	$\mathbb{E} [\ \mathbf{g}(\lambda, \eta)\ _*^2 \mathcal{F}_t] < G^2$	$\mathcal{O}\left(\frac{G^2 \log T}{\log \rho^{-1}\sqrt{T}}\right)$	Duchi <i>et al.</i> (2012) Corollary 3.5
SGD-Nesterov ²	$\gamma_t = 2/(t+1)$ $\beta_t = \frac{1}{2L\sqrt{t+1}}$	$\ \mathbf{g}(\lambda, \eta)\ _2 < G$	$\mathcal{O}\left(\frac{G^2 \log T}{\sqrt{T}}\right)$	Doan <i>et al.</i> (2020a) Theorem 2
SGD ³	$\gamma_t = \gamma/t$ $\gamma = \min\{1/2L, 2L/\mu\}$	$\ \mathbf{g}(\lambda, \eta)\ _* < G(\ \lambda\ _2 + 1)$	$\mathcal{O}\left(\frac{G^2 \log T}{T}\right)$	Doan <i>et al.</i> (2020b) Theorem 1,2

Notation: ¹ \mathcal{F}_t is the σ -field formed by all the iterates η_t, λ_t up to the t th MCGD iteration and $\|\mathbf{x}\|_*$ is the dual norm such that $\|\mathbf{x}\|_* = \sup_{\|\mathbf{z}\| \leq 1} \langle \mathbf{x}, \mathbf{z} \rangle$. ² β_t is the stepsize of the momentum. ³ L is the Lipschitz smoothness constant. ³ μ is the strong convexity constant.

Convergence and the Mixing Rate ρ Duchi *et al.* (2012) was the first to provide an analysis of the general MCGD setting. Their convergence rate is dependent on the mixing rate through the $1/\log \rho^{-1}$ term. For MCSA, this results in an overly conservative rate since, on challenging problems, mixing can be slow such that $\rho \approx 1$. Fortunately, Doan *et al.* (2020a,b) have recently shown that it is possible to obtain a rate independent of the mixing rate ρ . For example, in the result of Doan *et al.* (2020b), the influence of ρ decreases in a rate of $\mathcal{O}(1/T^2)$. This observation is critical since it implies that *trading “lower gradient variance” with a “slower mixing rate” could be profitable*. We exploit this observation in our novel MCSA scheme in Section 5.

Gradient Bound G Except for Doan *et al.* (2020b), most results assume that the gradient is bounded for $\forall \eta, \lambda$ as $\|\mathbf{g}(\lambda, \eta)\| < G$. Admittedly, this condition is strong, but it is similar to the bounded variance assumption $\mathbb{E}[\|\mathbf{g}\|^2] < G^2$ used in vanilla SGD, which is also known to be strong as it contradicts strong convexity (Nguyen *et al.*, 2018). Nonetheless, assuming G can have practical benefits beyond theoretical settings. For example, Geffner & Domke (2020) use G to compare the performance different VI gradient estimators. In a similar spirit, we will obtain the gradient bound G of different MCSA algorithms and compare their theoretical performance.

4 Demystifying Prior Markov Chain Score Ascent Methods

In this section, we will show that MSC and JSA both qualify as MCSA methods. Furthermore, we establish **(i)** the mixing rate of their implicitly defined kernel P and **(ii)** the upper bound on their gradient variance. This will provide insight into their practical non-asymptotic performance.

4.1 Technical Assumptions

To cast previous methods into MCSA, we need some technical assumptions.

Assumption 2. (*Bounded importance weight*) *The importance weight ratio $w(\mathbf{z}) = \pi(\mathbf{z})/q(\mathbf{z}; \lambda)$ is bounded by some finite constant as $w^* < \infty$ for all $\lambda \in \Lambda$ such that $\rho = (1 - 1/w^*) < 1$.*

This assumption is necessary to ensure Assumption 1, and can be practically ensured by using a variational family with heavy tails (Domke & Sheldon, 2018) or using a defensive mixture (Hesterberg, 1995; Holden *et al.*, 2009) as

$$q_{\text{def.}}(\mathbf{z}; \lambda) = w q(\mathbf{z}; \lambda) + (1 - w) \nu(\mathbf{z})$$

where $0 < w < 1$ and $\nu(\cdot)$ is a heavy tailed distribution such that $\sup_{\mathbf{z} \in \mathcal{Z}} \pi(\mathbf{z})/\nu(\mathbf{z}) < \infty$. Note that $q_{\text{def.}}$ is only used in the Markov chain kernels and $q(\cdot; \lambda^*)$ is still the output of the VI procedure. While these tricks help escape slowly mixing regions, this benefit quickly vanishes as we converge. Therefore, ensuring Assumption 2 seems unnecessary in practice unless we absolutely care about ergodicity (e.g. adaptive MCMC, Brofos *et al.* 2022; Holden *et al.* 2009).

Model (Variational Family) Misspecification and w^* Note that w^* is bounded below exponentially by the inclusive KL as shown in Proposition 2. Therefore, w^* will be large **(i)** in the initial steps of VI and **(ii)** under model (variational family) misspecification.

Assumption 3. (*Bounded Score*) *The score gradient is bounded for $\forall \lambda \in \Lambda$ and $\forall \mathbf{z} \in \mathcal{Z}$ such that $\|\mathbf{s}(\lambda; \mathbf{z})\|_* \leq L$ for some finite constant $L > 0$.*

Although this assumption is strong, it enables us to compare the gradient variance of MCSA methods. We empirically justify the bounds obtained using Assumption 3 in Section 6.2.

4.2 Markovian Score Climbing

MSC (Algorithm 3) is a simple instance of MCSA where $\eta_t = \mathbf{z}_t$ and $P_{\lambda_t} = K_{\lambda_t}$, where K_{λ_t} is the conditional importance sampling (CIS) kernel (originally proposed by Andrieu *et al.* (2018)) where the proposals are generated from $q(\cdot; \lambda_t)$. Although MSC uses only a single sample for the Markov chain, the CIS kernel internally uses N proposals to generate a single sample. Therefore, N in MSC has a different meaning, but it still indicates the computational budget.

Theorem 1. MSC (Naesseth *et al.*, 2020) is obtained by defining

$$P_{\lambda}^n(\eta, d\eta') = K_{\lambda}^n(\mathbf{z}, d\mathbf{z}')$$

with $\eta_t = \mathbf{z}_t$ where $K_{\lambda}(\mathbf{z}, \cdot)$ is the CIS kernel with $q_{\text{def.}}(\cdot; \lambda)$ as its proposal distribution. Then, given Assumption 2 and 3, the mixing rate and the gradient bounds are given as

$$d_{\text{TV}}(P_{\lambda}^n(\eta, \cdot), \Pi) \leq \left(1 - \frac{N-1}{2w^* + N-2}\right)^n \quad \text{and} \quad \mathbb{E}\left[\|\mathbf{g}(\lambda, \eta)\|_*^2 \mid \mathcal{F}_t\right] \leq L^2,$$

where $w^* = \sup_{\mathbf{z}} \pi(\mathbf{z}) / q_{\text{def.}}(\mathbf{z}; \lambda)$.

Proof. See page 18 of the *supplementary material*.

Discussion Theorem 1 shows that the gradient variance of MSA is insensitive to N . Although the mixing rate does improve with N , when w^* is large due to model misspecification and lack of convergence (see the discussion in Section 4.1), this will be marginal. Overall, *the performance of MSC cannot be improved by increasing the computational budget N* . Even under stationarity, the variance of η is equal to that of a *single* posterior sample at best.

4.3 Joint Stochastic Approximation

JSA (Algorithm 4) was proposed for deep generative models where the likelihood factorizes into each datapoint, for which subsampling can be used through a random-scan version of the independent Metropolis-Hastings (IMH, Hastings 1970) kernel. Instead, we consider the general version of JSA with the vanilla IMH kernel since it can be used with any type of likelihood. At each MCGD step, JSA performs multiple Markov chain transitions and estimates the gradient by averaging all the intermediate states, which is closer to traditional MCMC estimation.

Independent Metropolis-Hastings Similarly to MSC, the IMH kernel in JSA generates proposals from $q(\cdot; \lambda_t)$. To show the geometric ergodicity of the implicit kernel P , we utilize the geometric convergence rate of IMH kernels provided by Mengerson & Tweedie (1996, Theorem 2.1) and Wang (2020). Furthermore, to derive an upper bound on the gradient variance, we use the exact n -step marginal IMH kernel derived by Smith & Tierney (1996) as

$$K_{\lambda}^n(\mathbf{z}, d\mathbf{z}') = T_n(w(\mathbf{z}) \vee w(\mathbf{z}')) \pi(\mathbf{z}') d\mathbf{z}' + \lambda^n(w(\mathbf{z})) \delta_{\mathbf{z}}(d\mathbf{z}') \quad (2)$$

where $w(\mathbf{z}) = \pi(\mathbf{z}) / q_{\text{def.}}(\mathbf{z}; \lambda)$, $x \vee y = \max(x, y)$, and for $R(w) = \{\mathbf{z}' \mid w(\mathbf{z}') \leq w\}$,

$$T_n(w) = \int_w^{\infty} \frac{n}{v^2} \lambda^{n-1}(v) dv, \quad \text{and} \quad \lambda(w) = \int_{R(w)} \left(1 - \frac{w(\mathbf{z}')}{w}\right) \pi(d\mathbf{z}'). \quad (3)$$

Theorem 2. JSA (Ou & Song, 2020) is obtained by defining

$$P_{\lambda}^n(\eta, d\eta') = K_{\lambda}^{N(n-1)+1}(\mathbf{z}^{(1)}, d\mathbf{z}'^{(1)}) K_{\lambda}^{N(n-1)+2}(\mathbf{z}^{(2)}, d\mathbf{z}'^{(2)}) \cdots K_{\lambda}^{N(n-1)+N}(\mathbf{z}^{(N)}, d\mathbf{z}'^{(N)})$$

with $\eta_t = (\mathbf{z}_t^{(1)}, \mathbf{z}_t^{(2)}, \dots, \mathbf{z}_t^{(N)})$. Then, given Assumption 2 and 3, the mixing rate and the gradient variance bounds are

$$d_{\text{TV}}(P_{\lambda}^n(\eta, \cdot), \Pi) \leq C(\rho, N) \rho^{nN} \quad \text{and} \quad \mathbb{E}\left[\|\mathbf{g}(\lambda, \eta)\|_*^2 \mid \mathcal{F}_t\right] \leq L^2 \left[\frac{1}{2} + \frac{3}{2} \frac{1}{N} + \mathcal{O}(1/w^*) \right],$$

where $w^* = \sup_{\mathbf{z}} \pi(\mathbf{z}) / q_{\text{def.}}(\mathbf{z}; \lambda)$ and $C(\rho, N) > 0$ is a finite constant depending on ρ and N .

Proof. See page 20 of the *supplementary material*.

Discussion According to Theorem 2, JSA benefits from increasing N both in terms of faster mixing rate and lower gradient variance. However, similarly to MSC, under lack of convergence and model misspecification (large w^*), the improvement becomes marginal. More importantly, *in the large w^* regime, the variance reduction is limited by the constant 1/2 term*. This is because a large w^* results in more Metropolis-Hastings rejections, increasing the correlation between the N samples. This inefficiency persists even under stationarity unless $\pi(\cdot) \in \mathcal{Q}$ and $q(\cdot; \lambda) = \pi(\cdot)$.

5 Parallel Markov Chain Score Ascent

Our analysis in Section 4 suggests that the statistical performance of both MSC and JSA are heavily affected by model specification and the state of convergence through w^* . This affects both the mixing rate ρ and the gradient variance G . Furthermore, a large w^* abolishes our ability to counterbalance the inefficiency by increasing the computational budget N . However, ρ and G do not impact convergence equally. In particular, we noted in Section 3.1 that recent results on MCGD suggest that MCSA depends more on the gradient variance than the mixing rate. We turn to leverage this understanding to overcome the limitations of previous MCSA methods.

5.1 Parallel Markov Chain Score Ascent

We propose a novel scheme, *parallel Markov chain score ascent* (pMCSA, Algorithm 5), that embraces a slower mixing rate in order to consistently achieve an $\mathcal{O}(1/N)$ variance reduction, even on challenging problems with a large w^* ,

Algorithm Description Unlike JSA, that uses N sequential Markov chain states, pMCSA operates N parallel Markov chains. To maintain a similar per-iteration cost with JSA, it performs only a single Markov chain transition for each chain. Since the chains are independent, the Metropolis-Hastings rejections do not affect the variance of pMCSA.

Theorem 3. pMCSA, our proposed scheme, is obtained by setting

$$P_\lambda^n(\eta, d\eta') = K_\lambda^n(\mathbf{z}^{(1)}, d\mathbf{z}'^{(1)}) K_\lambda^n(\mathbf{z}^{(2)}, d\mathbf{z}'^{(2)}) \cdot \dots \cdot K_\lambda^n(\mathbf{z}^{(N)}, d\mathbf{z}'^{(N)})$$

with $\eta = (\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(N)})$. Then, given Assumption 2 and 3, the mixing rate and the gradient variance bounds are

$$d_{\text{TV}}(P_\lambda^n(\eta, \cdot), \Pi) \leq C(N) \rho^n \quad \text{and} \quad \mathbb{E}\left[\|\mathbf{g}(\lambda, \eta)\|_*^2 \mid \mathcal{F}_t\right] \leq L^2 \left[\frac{1}{N} + \frac{1}{N} \left(1 - \frac{1}{w^*}\right) \right],$$

where $w^* = \sup_{\mathbf{z}} \pi(\mathbf{z}) / q_{\text{def.}}(\mathbf{z}; \lambda)$ and C is some positive constant depending on N .

Proof. See page 22 of the *supplementary material*.

Discussion Unlike JSA and MSC, the variance reduction rate of pMCSA is independent of w^* . Therefore, it should perform significantly better on challenging practical problems. If we consider the rate of Duchi *et al.* (2012), the combined rate is constant with respect to N since it cancels out. In practice, however, we observe that increasing N accelerates convergence quite dramatically. Therefore, the mixing rate independent convergence rates by Doan *et al.* (2020a,b) appears to better reflect practical performance. This is because **(i)** the mixing rate ρ is a conservative *global* bound and **(ii)** the mixing rate will improve naturally as MCSA converges.

5.2 Computational Cost Comparison

The three schemes using the CIS and IMH kernels have different computational costs depending on N as organized in Table 2.

Cost of Sampling Proposals For the CIS kernel used by MSC, N controls the number of internal proposals sampled from $q(\cdot; \lambda)$. For JSA and pMCSA, the IMH kernel only uses a single sample from $q(\cdot; \lambda)$, but applies the kernel N times. On the other hand, pMCSA needs twice more evaluations of $q(\cdot; \lambda)$. However, this added cost is minimal since it is dominated by that of evaluating $p(\mathbf{z}, \mathbf{x})$.

Cost of Estimating the Score When estimating the score, MSC computes $\nabla_\lambda \log q(\mathbf{z}; \lambda)$ only once, while JSA and our proposed scheme compute it N times. However, Naesseth *et al.* (2020) also discuss a Rao-Blackwellized version of the CIS kernel, which also computes the score N times. Lastly, notice that MCSA methods do not need to differentiate through the likelihood $p(\mathbf{z}, \mathbf{x})$, unlike ELBO maximization, making its per-iteration cost significantly cheaper.

Table 2: Computational Costs

	Applying P_λ			Estimating \mathbf{g}	
	$p(\mathbf{z}, \mathbf{x})$	$q(\mathbf{z}; \lambda)$	$q(\mathbf{z}; \lambda)$	$p(\mathbf{z}, \mathbf{x})$	$q(\mathbf{z}; \lambda)$
	# Eval.	# Eval.	# Samples	# Grad.	# Grad.
ELBO	0	0	N	N	N
MSC	$N - 1$	N	$N - 1$	0	¹ Vanilla CIS kernel. ² Rao-Blackwellized CIS kernel.
JSA	N	$N + 1$	N	0	N
pMCSA	N	$2N$	N	0	N

Note: We assume that the variables are cached as much as possible. ¹Vanilla CIS kernel. ²Rao-Blackwellized CIS kernel.

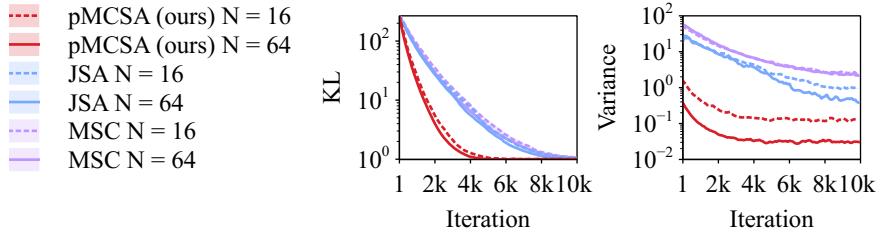


Figure 1: **KL divergence and gradient variance.** pMCSA not only achieves the least gradient variance, but its variance also scales better with N . The target distribution is a 20-D multivariate Gaussian with $\nu = 100$. The error bands are the 80% quantiles, while the solid lines are the median of 20 replications.

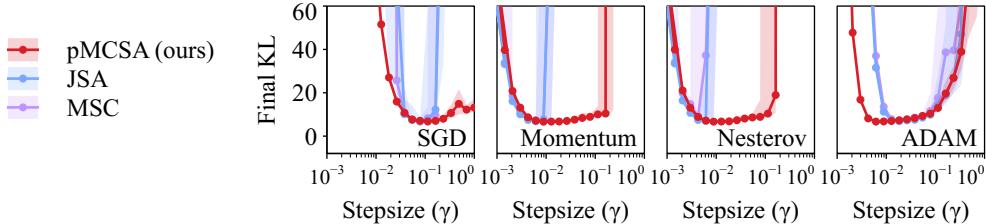


Figure 2: **Optimizer stepsize (γ) versus final KL.** pMCSA is the least sensitive to optimizer hyperparameters and results in stable convergence. The final KL is obtained at the 10^4 th iteration. The target distribution is a 100-D Gaussian with $\nu = 500$. The error bands are the 80% quantiles, while the solid lines are the median of 20 replications.

6 Evaluations

6.1 Experimental Setup

Implementation For the realistic experiments, we implemented MCSA methods on top of the Turing (Ge *et al.*, 2018) probabilistic programming framework². For the variational family, we use diagonal multivariate Gaussians with the support transformation of Kucukelbir *et al.* (2017). We use the ADAM optimizer by Kingma & Ba (2015) with a stepsize of 0.01 in all experiments. The budget is set to $N = 10$ for all experiments unless specified.

Baselines We compare **(i)** pMCSA (ours, Section 5), **(ii)** JSA (Ou & Song, 2020), **(iii)** MSC (Naesseth *et al.*, 2020), **(iv)** MSC with Rao-Blackwellization (**MSC-RB**, Naesseth *et al.* 2020), and **(v)** evidence lower-bound maximization (**ELBO**, Kucukelbir *et al.* 2017; Ranganath *et al.* 2014) with the path derivative estimator (Roeder *et al.*, 2017).

6.2 Simulations

Setup First, we verify our theoretical analysis on multivariate Gaussians with full-rank covariances sampled from Wishart distribution with ν degrees of freedom (values of ν are in the figure captions). This problem is challenging since an IMH (used by pMCSA, JSA) or CIS (used by MSC, MSC-RB) kernel with a diagonal Gaussian proposal will mix slowly due to a large w^* .

Gradient Variance We evaluate our theoretical analysis of the gradient variance. The variance is estimated from 512 independent Markov chains using the parameters generated by the main MCSA procedure. The estimated variances are shown in Figure 1. We make the following observations: **(i)** pMCSA has the lowest variance overall, and it consistently benefits from increasing N . **(ii)** MSC does not benefit from increasing N whatsoever. **(iii)** JSA does not benefit from increasing N until $q(\cdot; \lambda)$ has sufficiently converged (when w^* has become small). These results confirm our theoretical results in Sections 4 and 5.

Robustness Against Optimizers Since the convergence of most sophisticated SGD optimizers has yet to be established for MCGD, we empirically investigate their effectiveness. The results using SGD (Bottou *et al.*, 2018; Robbins & Monroe, 1951), Momentum (Polyak, 1964), Nesterov (Nesterov, 1983), ADAM (Kingma & Ba, 2015), and varying stepsizes are shown in Figure 2. Clearly, pMCSA successfully converges for the broadest variety of optimizer settings. Overall, most MCSA methods seem to be the most stable with ADAM, which points out that establishing the convergence of ADAM for MCGD will be a promising direction for future works.

²Available at <https://github.com/Red-Portal/KLpqVI.jl>

Table 3: Test Log Predictive Density on **Bayesian Neural Network Regression**

	D_λ	D_x	N_{train}	ELBO		MCSA Variants		
				$N = 1$	$N = 10$	pMCSA (ours)	JSA	MSC
yacht	403	6	277	-2.45 ± 0.01	-2.44 ± 0.01	-2.49 ± 0.01	-3.00 ± 0.05	-2.98 ± 0.04
concrete	503	8	927	-3.25 ± 0.01	-3.24 ± 0.01	-3.20 ± 0.01	-3.33 ± 0.02	-3.32 ± 0.02
airfoil	353	6	1352	-2.53 ± 0.02	-2.56 ± 0.02	-2.27 ± 0.02	-2.51 ± 0.02	-2.53 ± 0.01
energy	503	9	691	-2.42 ± 0.02	-2.40 ± 0.02	-1.92 ± 0.03	-2.38 ± 0.02	-2.37 ± 0.02
wine	653	12	1439	-0.96 ± 0.01	-0.96 ± 0.01	-0.95 ± 0.01	-0.97 ± 0.01	-0.97 ± 0.01
boston	753	14	455	-2.72 ± 0.03	-2.70 ± 0.03	-2.69 ± 0.02	-2.82 ± 0.02	-2.80 ± 0.03
sml	1203	23	3723	-1.32 ± 0.01	-1.25 ± 0.02	-1.22 ± 0.01	-1.72 ± 0.01	-1.97 ± 0.02
gas	6503	129	2308	-0.06 ± 0.01	0.13 ± 0.03	-0.09 ± 0.02	-0.47 ± 0.03	-0.47 ± 0.04

Table 4: Test Log Predictive Density on **Robust Gaussian Process Regression**

	D_λ	D_x	N_{train}	ELBO		MCSA Variants		
				$N = 1$	pMCSA (ours)	JSA	MSC	MSC-RB
yacht	287	6	277	-3.63 ± 0.02	-3.31 ± 0.04	-3.29 ± 0.05	-3.25 ± 0.04	-3.27 ± 0.05
airfoil	353	6	1352	-3.14 ± 0.01	-2.63 ± 0.01	-2.83 ± 0.04	-2.77 ± 0.02	-2.73 ± 0.02
boston	472	13	455	-2.98 ± 0.01	-2.96 ± 0.02	-3.00 ± 0.03	-3.00 ± 0.03	-2.96 ± 0.03
energy	703	8	691	-2.75 ± 0.01	-2.58 ± 0.03	-2.78 ± 0.04	-2.70 ± 0.04	-2.72 ± 0.05
concrete	939	8	927	-3.68 ± 0.01	-3.49 ± 0.01	-3.69 ± 0.02	-3.59 ± 0.04	-3.57 ± 0.02
wine	1454	11	1439	-1.02 ± 0.01	-0.94 ± 0.02	-1.04 ± 0.01	-1.00 ± 0.02	-0.99 ± 0.02
gas	2440	128	2308	0.18 ± 0.02	-0.86 ± 0.02	-1.10 ± 0.03	-1.10 ± 0.04	-1.06 ± 0.02

¹ D_λ : Dimensionality of λ , D_x : Number of features, N_{train} : Number of training data points.

² \pm denotes the 95% bootstrap confidence intervals obtained from 20 replications.

³ Bolded numbers don't have enough evidence to be distinguished from the best performing method under a .05 significance threshold (Friedman test with Nemenyi post-hoc test, Demšar 2006).

6.3 Bayesian Neural Network Regression

Setup For realistic experiments, we train Bayesian neural networks (BNN, Neal 1996) for regression. We use datasets from the UCI repository (Dua & Graff, 2017) with 90% random train-test splits. We run all methods with $T = 5 \cdot 10^4$ iterations. For the model, we use the priors and forward propagation method of Hernandez-Lobato & Adams (2015) with a 50-unit hidden layer (see Appendix C.1).

Results The results are shown in Table 3. pMCSA achieves the best performance compared to all other MCSA methods. Also, its overall performance is comparable to exclusive KL minimization methods (ELBO) unlike other MCSA methods. Furthermore, on airfoil and energy, pMCSA improves over ELBO by 0.29 nat and 0.48 nat. Even on gas where pMCSA did not beat ELBO, its performance is comparable, and it dominates all other MCSA methods by roughly 0.4 nat. Additional experimental results can be found in Appendix E.1

6.4 Robust Gaussian Process Regression

Setup We train Gaussian processes (GP) with a Student-T likelihood for robust regression. We use datasets from the UCI repository (Dua & Graff, 2017) with 90% random train-test splits. We use the Matérn 5/2 covariance kernel with automatic relevance determination (Neal, 1996) (see Appendix C.2). We run all methods with $T = 2 \cdot 10^4$ iterations. For prediction, we use the mode of $q(\cdot; \lambda)$ for the hyperparameters and marginalize the latent function over $q(\cdot; \lambda)$ (Rasmussen & Williams, 2006). We consider ELBO with only $N = 1$ since differentiating through the likelihood makes its per-iteration cost comparable to MCSA methods with $N = 10$.

Results The results are shown in Table 4. Except for gas, pMCSA achieves better performance than all other methods. This suggests that, overall, the exclusive KL may be less effective for GP posteriors. Although ELBO achieves the best performance on gas, pMCSA dominates other MCSA methods. Our encouraging regression results suggest that incorporating methods such as inducing points (Snelson & Ghahramani, 2005) into MCSA may lead to an important new class of GP models. Additional experimental results can be found in Appendix E.2.

7 Related Works

Inclusive KL minimization Our MCSA framework generalizes MSC (Naesseth *et al.*, 2020) and JSA Ou & Song (2020), which are inclusive KL minimization based on SGD and Markov chains. Similar to MCSA is the method of Li *et al.* (2017). However, the convergence of this method is not guaranteed since it uses short Markov chains, disqualifying for MCSA. Other methods based on biased gradients have been proposed by Bornschein & Bengio (2015); Le *et al.* (2019), but these are specific for deep generative models. On a different note, Jerfel *et al.* (2021) use boosting instead of SGD to minimize the inclusive KL, which gradually builds a complex variational approximation from a simple variational family.

Beyond the KL Divergence Discovering alternative divergences for VI has been an active research area. For example, the χ^2 (Dieng *et al.*, 2017), f (Wan *et al.*, 2020; Wang *et al.*, 2018), α (Hernandez-Lobato *et al.*, 2016; Li & Turner, 2016; Regli & Silva, 2018) divergences have been studied for VI. However, for gradient estimation these methods involve the importance ratio $w(\mathbf{z}) = \pi(\mathbf{z})/q(\mathbf{z})$, which leads to significant variance and low signal-to-noise ratio under model misspecification (Geffner & Domke, 2021a,c). In contrast, under stationarity, the variance of pMCSA is σ^2/N (σ is the variance of the score over the posterior) regardless of model misspecification. Meanwhile, Geffner & Domke (2021b); Ruiz & Titsias (2019); Salimans *et al.* (2015); Zhang *et al.* (2021) construct implicit divergences formed by MCMC. It is to be seen how these divergences compare against the inclusive KL divergence.

Adaptive MCMC and MCSA As pointed out by Ou & Song (2020), using $q(\cdot; \lambda)$ within the MCMC kernel makes MCSA structurally equivalent to adaptive MCMC. In particular, Andrieu & Thoms (2008); Brofos *et al.* (2022); Gabrié *et al.* (2022); Garthwaite *et al.* (2016) discuss the use of stochastic approximation in adaptive MCMC. Also, Andrieu & Moulines (2006); Brofos *et al.* (2022); Giordani & Kohn (2010); Habib & Barber (2019); Holden *et al.* (2009); Keith *et al.* (2008); Neklyudov *et al.* (2019) specifically discuss adapting the proposal of IMH kernels, and some of them use KL divergence minimization. These methods focus on showing ergodicity the samples $(\mathbf{z}^{(n)})$ in our notation not the convergence of the variational approximation (or proposal in an adaptive MCMC context) $q(\cdot; \lambda)$. In this work, we focused on the convergence of $q(\cdot; \lambda)$, which could advance the adaptive MCMC side of the story.

8 Discussions

This paper presented a new theoretical framework for analyzing inclusive KL divergence minimization methods based on running SGD with Markov chains. Furthermore, we proposed pMCSA, a new MCSA method that enjoys substantially low variance. We have shown that this theoretical improvement translates into better empirical performance.

Limitations Our work has two main limitations. Firstly, since our work aims to understand existing MCSA methods, it inherits the current limitations of MCSA methods. For example, minibatch subsampling is challenging for models with non-factorizable likelihoods (Naesseth *et al.*, 2020). Secondly, our theoretical analysis in Section 4 requires Assumption 3, which is strong, but required to connect with MCGD. An important future direction would be to relax the assumptions needed by MCGD.

Towards Alternative Divergences In Section 6, we have shown that minimizing the *inclusive* KL is competitive against minimizing the *exclusive* KL on general Bayesian inference problems. Although Dhaka *et al.* (2021) has shown that the inclusive KL fails on high-dimensional problems, this is only the case under the presence of strong correlations. Before entirely ditching the inclusive KL, it is essential to ask, “how correlated posteriors really are in practice?” Furthermore, the true performance of alternative divergences is often masked by the limitations of the inference procedure (Geffner & Domke, 2021a,c). Given that pMCSA significantly advances the best-known performance of inclusive KL minimization, it is possible that similar improvements could be extracted from other divergences. To conclude, our results motivate further development of better inference algorithms for alternative divergence measures.

Acknowledgments and Disclosure of Funding

This work initially started in the process of understanding the performance of Markovian score climbing. We thank Hongseok Yang for pointing us to a relevant related work, Guanyang Wang

for insightful discussions about the independent Metropolis-Hastings algorithm, Geon Park and Kwanghee Choi for constructive comments that enriched this paper. We also aknowledge the Computer Science Department of Sogang University for providing computational resources.

References

- Andrieu, Christophe, & Moulines, Éric. 2006. On the Ergodicity Properties of Some Adaptive MCMC Algorithms. *The Annals of Applied Probability*, **16**(3).
- Andrieu, Christophe, & Thoms, Johannes. 2008. A Tutorial on Adaptive MCMC. *Statistics and Computing*, **18**(4), 343–373.
- Andrieu, Christophe, Lee, Anthony, & Vihola, Matti. 2018. Uniform Ergodicity of the Iterated Conditional SMC and Geometric Ergodicity of Particle Gibbs Samplers. *Bernoulli*, **24**(2).
- Blei, David M., Kucukelbir, Alp, & McAuliffe, Jon D. 2017. Variational Inference: A Review for Statisticians. *Journal of the American Statistical Association*, **112**(518), 859–877.
- Bornschein, Jörg, & Bengio, Yoshua. 2015 (May). Reweighted Wake-Sleep. In: *Proceedings of the International Conference on Learning Representations*.
- Bottou, Léon, Curtis, Frank E., & Nocedal, Jorge. 2018. Optimization Methods for Large-Scale Machine Learning. *SIAM Review*, **60**(2), 223–311.
- Brofos, James, Gabrie, Marylou, Brubaker, Marcus A., & Lederman, Roy R. 2022. Adaptation of the Independent Metropolis-Hastings Sampler with Normalizing Flow Proposals. Pages 5949–5986 of: *Proceedings of the International Conference on Artificial Intelligence and Statistics*. PMLR, vol. 151. ML Research Press.
- de Boer, Pieter-Tjerk, Kroese, Dirk P., Mannor, Shie, & Rubinstein, Reuven Y. 2005. A Tutorial on the Cross-Entropy Method. *Annals of Operations Research*, **134**(1), 19–67.
- Debavelaere, Vianney, Durrelman, Stanley, & Alloussonnière, Stéphanie. 2021. On the Convergence of Stochastic Approximations under a Subgeometric Ergodic Markov Dynamic. *Electronic Journal of Statistics*, **15**(1).
- Demšar, Janez. 2006. Statistical Comparisons of Classifiers over Multiple Data Sets. *Journal of Machine Learning Research*, **7**(1), 1–30.
- Dhaka, Akash Kumar, Catalina, Alejandro, Welandawe, Manushi, Andersen, Michael R., Huggins, Jonathan, & Vehtari, Aki. 2021. Challenges and Opportunities in High Dimensional Variational Inference. In: *Advances in Neural Information Processing Systems*.
- Dieng, Adji Bousso, Tran, Dustin, Ranganath, Rajesh, Paisley, John, & Blei, David. 2017. Variational Inference via χ Upper Bound Minimization. Pages 2729–2738 of: *Advances in Neural Information Processing Systems*, vol. 30. Long Beach, California, USA: Curran Associates, Inc.
- Doan, Thinh T., Nguyen, Lam M., Pham, Nhan H., & Romberg, Justin. 2020a (Oct.). *Convergence Rates of Accelerated Markov Gradient Descent with Applications in Reinforcement Learning*. Tech. rept. arXiv:2002.02873 [math]. ArXiv.
- Doan, Thinh T., Nguyen, Lam M., Pham, Nhan H., & Romberg, Justin. 2020b (Apr.). *Finite-Time Analysis of Stochastic Gradient Descent under Markov Randomness*. Tech. rept. arXiv:2003.10973. ArXiv.
- Domke, Justin, & Sheldon, Daniel R. 2018. Importance Weighting and Variational Inference. In: *Advances in Neural Information Processing Systems*, vol. 31. Curran Associates, Inc.
- Dua, Dheeru, & Graff, Casey. 2017. UCI Machine Learning Repository.
- Duchi, John C., Agarwal, Alekh, Johansson, Mikael, & Jordan, Michael I. 2012. Ergodic Mirror Descent. *SIAM Journal on Optimization*, **22**(4), 1549–1578.
- Gabrié, Marylou, Rotskoff, Grant M., & Vanden-Eijnden, Eric. 2022. Adaptive Monte Carlo Augmented with Normalizing Flows. *Proceedings of the National Academy of Sciences*, **119**(10), e2109420119.
- Garthwaite, P. H., Fan, Y., & Sisson, S. A. 2016. Adaptive Optimal Scaling of Metropolis–Hastings Algorithms Using the Robbins–Monro Process. *Communications in Statistics - Theory and Methods*, **45**(17), 5098–5111.

- Ge, Hong, Xu, Kai, & Ghahramani, Zoubin. 2018. Turing: A Language for Flexible Probabilistic Inference. *Pages 1682–1690 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 84. ML Research Press.
- Geffner, Tomas, & Domke, Justin. 2020. A Rule for Gradient Estimator Selection, with an Application to Variational Inference. *Pages 1803–1812 of: Proceedings of the International Conference on Artificial Intelligence and Statistics*. PMLR, vol. 108. ML Research Press.
- Geffner, Tomas, & Domke, Justin. 2021a. Empirical Evaluation of Biased Methods for Alpha Divergence Minimization. *In: Proceedings of the Symposium on Advances in Approximate Bayesian Inference*.
- Geffner, Tomas, & Domke, Justin. 2021b. MCMC Variational Inference via Uncorrected Hamiltonian Annealing. *Pages 639–651 of: Advances in Neural Information Processing Systems*, vol. 34. Curran Associates, Inc.
- Geffner, Tomas, & Domke, Justin. 2021c. On the Difficulty of Unbiased Alpha Divergence Minimization. *Pages 3650–3659 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 139. ML Research Press.
- Giordani, Paolo, & Kohn, Robert. 2010. Adaptive Independent Metropolis–Hastings by Fast Estimation of Mixtures of Normals. *Journal of Computational and Graphical Statistics*, **19**(2), 243–259.
- Gu, Ming Gao, & Kong, Fan Hui. 1998. A Stochastic Approximation Algorithm with Markov Chain Monte-Carlo Method for Incomplete Data Estimation Problems. *Proceedings of the National Academy of Sciences*, **95**(13), 7270–7274.
- Habib, Raza, & Barber, David. 2019. Auxiliary Variational MCMC. *In: Proceedings of the International Conference on Learning Representations*.
- Hastings, W. K. 1970. Monte Carlo Sampling Methods Using Markov Chains and Their Applications. *Biometrika*, **57**(1), 97–109.
- Hernandez-Lobato, Jose, Li, Yingzhen, Rowland, Mark, Bui, Thang, Hernandez-Lobato, Daniel, & Turner, Richard. 2016. Black-Box Alpha Divergence Minimization. *Pages 1511–1520 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 48. New York, New York, USA: ML Research Press.
- Hernandez-Lobato, Jose Miguel, & Adams, Ryan. 2015. Probabilistic Backpropagation for Scalable Learning of Bayesian Neural Networks. *Pages 1861–1869 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 37. Lille, France: ML Research Press.
- Hesterberg, Tim. 1995. Weighted Average Importance Sampling and Defensive Mixture Distributions. *Technometrics*, **37**(2), 185–194.
- Holden, Lars, Hauge, Ragnar, & Holden, Marit. 2009. Adaptive Independent Metropolis–Hastings. *The Annals of Applied Probability*, **19**(1).
- Jerfel, Ghassen, Wang, Serena, Wong-Fannjiang, Clara, Heller, Katherine A., Ma, Yian, & Jordan, Michael I. 2021. Variational Refinement for Importance Sampling Using the Forward Kullback-Leibler Divergence. *Pages 1819–1829 of: Proceedings of the International Conference on Uncertainty in Artificial Intelligence*. PMLR, vol. 161. ML Research Press.
- Karimi, Belhal, Miasojedow, Blazej, Moulines, Eric, & Wai, Hoi-To. 2019. Non-Asymptotic Analysis of Biased Stochastic Approximation Scheme. *Pages 1944–1974 of: Proceedings of the Annual Conference on Learning Theory*. PMLR, vol. 99. ML Research Press.
- Keith, Jonathan M., Kroese, Dirk P., & Sofronov, George Y. 2008. Adaptive Independence Samplers. *Statistics and Computing*, **18**(4), 409–420.
- Kingma, Diederik P., & Ba, Jimmy. 2015. Adam: A Method for Stochastic Optimization. *In: Proceedings of the International Conference on Learning Representations*.
- Kucukelbir, Alp, Tran, Dustin, Ranganath, Rajesh, Gelman, Andrew, & Blei, David M. 2017. Automatic Differentiation Variational Inference. *Journal of Machine Learning Research*, **18**(14), 1–45.
- Le, Tuan Anh, Kosirok, Adam R., Siddharth, N., Teh, Yee Whye, & Wood, Frank. 2019 (July). Revisiting Reweighted Wake-Sleep for Models with Stochastic Control Flow. *In: Proceedings of the International Conference on Uncertainty in Artificial Intelligence*.

- Li, Yingzhen, & Turner, Richard E. 2016. Rényi Divergence Variational Inference. *In: Advances in Neural Information Processing Systems*, vol. 29. Curran Associates, Inc.
- Li, Yingzhen, Turner, Richard E., & Liu, Qiang. 2017 (May). *Approximate Inference with Amortised MCMC*. Tech. rept. arXiv:1702.08343 [cs, stat]. ArXiv.
- MacKay, David J.C. 2001 (June). *Local Minima, Symmetry-Breaking, and Model Pruning in Variational Free Energy Minimization*. Technical Report.
- Mengersen, K. L., & Tweedie, R. L. 1996. Rates of Convergence of the Hastings and Metropolis Algorithms. *The Annals of Statistics*, **24**(1), 101–121.
- Minka, Tom. 2005 (Jan.). *Divergence Measures and Message Passing*. Tech. rept. MSR-TR-2005-173. Microsoft Research.
- Naesseth, Christian, Lindsten, Fredrik, & Blei, David. 2020. Markovian Score Climbing: Variational Inference with $\text{KL}(p||q)$. *Pages 15499–15510 of: Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc.
- Neal, Radford M. 1996. *Bayesian Learning for Neural Networks*. Lecture Notes in Statistics, vol. 118. New York, NY: Springer New York.
- Neklyudov, Kirill, Egorov, Evgenii, Shvchikov, Pavel, & Vetrov, Dmitry. 2019 (June). *Metropolis-Hastings View on Variational Inference and Adversarial Training*. Tech. rept. arXiv:1810.07151 [cs, stat]. ArXiv.
- Nesterov, Yurii Evgen'evich. 1983. A Method of Solving a Convex Programming Problem with Convergence Rate $\$O(\sqrt{\frac{1}{k^2}})$. *Doklady Akademii Nauk SSSR*, **269**(3), 543–547.
- Nguyen, Lam, Nguyen, Phuong Ha, van Dijk, Marten, Richtarik, Peter, Scheinberg, Katya, & Takac, Martin. 2018. SGD and Hogwild! Convergence without the Bounded Gradients Assumption. *Pages 3750–3758 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 80. ML Research Press.
- Ou, Zhijian, & Song, Yunfu. 2020. Joint Stochastic Approximation and Its Application to Learning Discrete Latent Variable Models. *Pages 929–938 of: Proceedings of the International Conference on Uncertainty in Artificial Intelligence*. PMLR, vol. 124. ML Research Press.
- Polyak, B.T. 1964. Some Methods of Speeding up the Convergence of Iteration Methods. *USSR Computational Mathematics and Mathematical Physics*, **4**(5), 1–17.
- Ram, S. Sundhar, Nedić, A., & Veeravalli, V. V. 2009. Incremental Stochastic Subgradient Algorithms for Convex Optimization. *SIAM Journal on Optimization*, **20**(2), 691–717.
- Ranganath, Rajesh, Gerrish, Sean, & Blei, David. 2014. Black Box Variational Inference. *Pages 814–822 of: Proceedings of the International Conference on Artificial Intelligence and Statistics*. PMLR, vol. 33. Reykjavik, Iceland: ML Research Press.
- Rasmussen, Carl Edward, & Williams, Christopher K. I. 2006. *Gaussian Processes for Machine Learning*. Adaptive Comput. Mach. Learn. Cambridge, Mass: MIT Press.
- Regli, Jean-Baptiste, & Silva, Ricardo. 2018. Alpha-Beta Divergence for Variational Inference. May.
- Robbins, Herbert, & Monro, Sutton. 1951. A Stochastic Approximation Method. *The Annals of Mathematical Statistics*, **22**(3), 400–407.
- Robert, Christian P., & Casella, George. 2004. *Monte Carlo Statistical Methods*. Springer Texts in Statistics. New York, NY: Springer New York.
- Roeder, Geoffrey, Wu, Yuhuai, & Duvenaud, David K. 2017. Sticking the Landing: Simple, Lower-Variance Gradient Estimators for Variational Inference. *In: Advances in Neural Information Processing Systems*, vol. 30. Curran Associates, Inc.
- Ruiz, Francisco, & Titsias, Michalis. 2019. A Contrastive Divergence for Combining Variational Inference and MCMC. *Pages 5537–5545 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 97. ML Research Press.
- Salimans, Tim, Kingma, Diederik, & Welling, Max. 2015. Markov Chain Monte Carlo and Variational Inference: Bridging the Gap. *Pages 1218–1226 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 37. Lille, France: ML Research Press.

- Smith, Richard L., & Tierney, Luke. 1996. *Exact Transition Probabilities for the Independence Metropolis Sampler*. Tech. rept.
- Snelson, Edward, & Ghahramani, Zoubin. 2005. Sparse Gaussian Processes Using Pseudo-Inputs. In: *Advances in Neural Information Processing Systems*, vol. 18. MIT Press.
- Sun, Tao, Sun, Yuejiao, & Yin, Wotao. 2018. On Markov Chain Gradient Descent. In: *Advances in Neural Information Processing Systems*, vol. 31. Curran Associates, Inc.
- Tadić, Vladislav B., & Doucet, Arnaud. 2017. Asymptotic Bias of Stochastic Gradient Search. *The Annals of Applied Probability*, 27(6).
- Trippé, Brian, & Turner, Richard. 2017. Overpruning in Variational Bayesian Neural Networks. Tech. rept. arXiv:1801.06230. ArXiv.
- Wan, Neng, Li, Dapeng, & Hovakimyan, Naira. 2020. F-Divergence Variational Inference. Pages 17370–17379 of: *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc.
- Wang, Dilin, Liu, Hao, & Liu, Qiang. 2018. Variational Inference with Tail-Adaptive f-Divergence. In: *Advances in Neural Information Processing Systems*, vol. 31. Curran Associates, Inc.
- Wang, Guanyang. 2020. Exact Convergence Rate Analysis of the Independent Metropolis-Hastings Algorithms. *To appear in Bernoulli*, Dec.
- Xiong, Huaqing, Xu, Tengyu, Liang, Yingbin, & Zhang, Wei. 2021. Non-Asymptotic Convergence of Adam-Type Reinforcement Learning Algorithms under Markovian Sampling. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(12), 10460–10468.
- Zhang, Guodong, Hsu, Kyle, Li, Jianing, Finn, Chelsea, & Grosse, Roger B. 2021. Differentiable Annealed Importance Sampling and the Perils of Gradient Noise. Pages 19398–19410 of: *Advances in Neural Information Processing Systems*, vol. 34. Curran Associates, Inc.

A Computational Resources

Table 5: Computational Resources for Bayesian Neural Network Regression

Type	Model and Specifications
System Topology	4 nodes with 20 logical threads each
Processor	Intel Xeon Xeon E5-2640 v4, 2.2 GHz (maximum 3.1 GHz)
Cache	32 kB L1, 256 kB L2, and 25 MB L3
Memory	64GB RAM

Table 6: Computational Resources for Robust Gaussian Process Regression

Type	Model and Specifications
System Topology	1 node with 16 logical threads
Processor	AMD EPYC 7262, 3.2 GHz (maximum 3.4 GHz)
Accelerator	NVIDIA Titan RTX, 1.3 GHZ, 24GB RAM
Cache	256 kB L1, 4MiB L2, and 128MiB L3
Memory	126GB RAM

B Pseudocodes

B.1 Markov Chain Kernels

Algorithm 1: Conditional Importance Sampling Kernel

Input: previous sample \mathbf{z}_{t-1} , previous parameter λ_{t-1} , number of proposals N
 $\mathbf{z}^{(0)} = \mathbf{z}_{t-1}$
 $\mathbf{z}^{(i)} \sim q_{\text{def.}}(\mathbf{z}; \lambda_{t-1}) \quad \text{for } i = 1, 2, \dots, N$
 $\tilde{w}(\mathbf{z}^{(i)}) = p(\mathbf{z}^{(i)}, \mathbf{x}) / q_{\text{def.}}(\mathbf{z}^{(i)}; \lambda_{t-1}) \quad \text{for } i = 0, 1, \dots, N$
 $\bar{w}^{(i)} = \frac{\tilde{w}(\mathbf{z}^{(i)})}{\sum_{i=0}^N \tilde{w}(\mathbf{z}^{(i)})} \quad \text{for } i = 0, 1, \dots, N$
 $\mathbf{z}_t \sim \text{Multinomial}(\bar{w}^{(0)}, \bar{w}^{(1)}, \dots, \bar{w}^{(N)})$

Algorithm 2: Independent Metropolis-Hastings Kernel

Input: previous sample \mathbf{z}_{t-1} , previous parameter λ_{t-1} , $\mathbf{z}^* \sim q_{\text{def.}}(\mathbf{z}; \lambda_{t-1})$
 $w(\mathbf{z}) = p(\mathbf{z}, \mathbf{x}) / q_{\text{def.}}(\mathbf{z}; \lambda_{t-1})$
 $\alpha = \min(w(\mathbf{z}^*) / w(\mathbf{z}_{t-1}), 1)$
 $u \sim \text{Uniform}(0, 1)$
if $u < \alpha$ **then**
 | $\mathbf{z}_t = \mathbf{z}^*$
else
 | $\mathbf{z}_t = \mathbf{z}_{t-1}$
end

B.2 Markov Chain Score Ascent Algorithms

Algorithm 3: Markovian Score Climbing

Input: MCMC kernel $K_\lambda(\mathbf{z}, \cdot)$, initial sample \mathbf{z}_0 ,
initial parameter λ_0 , number of iterations T ,
stepsize schedule γ_t

for $t = 1, 2, \dots, T$ **do**

$\mathbf{z}_t \sim K_{\lambda_{t-1}}(\mathbf{z}_{t-1}, \cdot)$

$\mathbf{g}(\lambda) = \mathbf{s}(\lambda; \mathbf{z}_t)$

$\lambda_t = \lambda_{t-1} + \gamma_t \mathbf{g}(\lambda_{t-1})$

end

Algorithm 4: Joint Stochastic Approximation

Input: MCMC kernel $K_\lambda(\mathbf{z}, \cdot)$, initial sample $\mathbf{z}_0^{(N)}$,
initial parameter λ_0 , number of iterations T ,
stepsize schedule γ_t

for $t = 1, 2, \dots, T$ **do**

$\mathbf{z}_t^{(0)} = \mathbf{z}_{t-1}^{(N)}$

for $n = 1, 2, \dots, N$ **do**

$\mathbf{z}_t^{(n)} \sim K_{\lambda_{t-1}}(\mathbf{z}_t^{(n-1)}, \cdot)$

end

$\mathbf{g}(\lambda) = \frac{1}{N} \sum_{n=1}^N \mathbf{s}(\lambda; \mathbf{z}_t^{(n)})$

$\lambda_t = \lambda_{t-1} + \gamma_t \mathbf{g}(\lambda_{t-1})$

end

Algorithm 5: Parallel Markov Chain Score Ascent

Input: MCMC kernel $K_\lambda(\mathbf{z}, \cdot)$, initial samples $\mathbf{z}_0^{(1)}, \dots, \mathbf{z}_0^{(N)}$,
initial parameter λ_0 , number of iterations T , stepsize
schedule γ_t

for $t = 1, 2, \dots, T$ **do**

for $n = 1, 2, \dots, N$ **do**

$\mathbf{z}_t^{(n)} \sim K_{\lambda_{t-1}}(\mathbf{z}_{t-1}^{(n)}, \cdot)$

end

$\mathbf{g}(\lambda) = \frac{1}{N} \sum_{i=1}^N \mathbf{s}(\lambda; \mathbf{z}_t^{(i)})$

$\lambda_t = \lambda_{t-1} + \gamma_t \mathbf{g}(\lambda_{t-1})$

end

C Probabilistic Models Used in the Experiments

C.1 Bayesian Neural Network Regression

We use the BNN model of Hernandez-Lobato & Adams (2015) defined as

$$\begin{aligned}\lambda^{-1} &\sim \text{Inverse-Gamma}(\alpha = 6, \beta = 6) \\ \gamma^{-1} &\sim \text{Inverse-Gamma}(\alpha = 6, \beta = 6) \\ \mathbf{W}_1 &\sim \mathcal{N}(\mathbf{0}, \lambda^{-1} \mathbf{I}) \\ \mathbf{z} &= \text{ReLU}(\mathbf{W}_1 \mathbf{x}_i) \\ \mathbf{W}_2 &\sim \mathcal{N}(\mathbf{0}, \lambda^{-1} \mathbf{I}) \\ \hat{y} &= \text{ReLU}(\mathbf{W}_2 \mathbf{z}) \\ y_i &\sim \mathcal{N}(\hat{y}, \gamma^{-1})\end{aligned}$$

where \mathbf{x}_i and y_i are the feature vector and target value of the i th datapoint. Given the variational distribution of $\lambda^{-1}, \gamma^{-1}, \mathbf{W}_1, \mathbf{W}_2$, we use the same posterior predictive approximation of Hernandez-Lobato & Adams (2015). We apply z-standardization (whitening) to the features \mathbf{x}_i and the target values y_i , and unwhiten the predictive distribution.

C.2 Robust Gaussian Process Logistic Regression

We perform robust Gaussian process regression by using a Student-T prior with a latent Gaussian process prior. The model is defined as

$$\begin{aligned} \log \sigma_f &\sim \mathcal{N}(0, 4) \\ \log \epsilon &\sim \mathcal{N}(0, 4) \\ \log \ell_i &\sim \mathcal{N}(0, 0.2) \\ f &\sim \mathcal{GP}\left(\mathbf{0}, \Sigma_{\sigma_f, \ell} + (\delta + \epsilon^2) \mathbf{I}\right) \\ \nu &\sim \text{Gamma}(\alpha = 4, \beta = 1/10) \\ \log \sigma_y &\sim \mathcal{N}(0, 4) \\ y_i &\sim \text{Student-T}(f(\mathbf{x}_i), \sigma_y, \nu). \end{aligned}$$

The covariance Σ is computed using a kernel $k(\cdot, \cdot)$ such that $[\Sigma]_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$ where \mathbf{x}_i and \mathbf{x}_j are data points in the dataset. For the kernel, we use the Matern 5/2 kernel with automatic relevance determination (Neal, 1996) defined as

$$k(\mathbf{x}, \mathbf{x}'; \sigma^2, \ell_1^2, \dots, \ell_D^2) = \sigma_f \left(1 + \sqrt{5}r + \frac{5}{3}r^2 \right) \exp(-\sqrt{5}r) \quad \text{where } r = \sum_{i=1}^D \frac{(\mathbf{x}_i - \mathbf{x}'_i)^2}{\ell_i^2}$$

where D is the number of dimensions. The jitter term δ is used for numerical stability. We set a small value of $\delta = 1 \times 10^{-6}$.

D Proofs

Proposition 1. Let $\eta = (\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(N)})$ and a Markov chain kernel $P_\lambda(\eta, \cdot)$ be Π -invariant where Π is defined as

$$\Pi(\eta) = \pi(\mathbf{z}^{(1)}) \pi(\mathbf{z}^{(2)}) \times \dots \times \pi(\mathbf{z}^{(N)}).$$

Then, by defining the target function f and the gradient estimator \mathbf{g} to be

$$f(\lambda, \eta) = -\frac{1}{N} \sum_{n=1}^N \log q(\mathbf{z}^{(n)}; \lambda) - \mathbb{H}[\pi], \quad \text{and} \quad \mathbf{g}(\lambda, \eta) = -\frac{1}{N} \sum_{n=1}^N \mathbf{s}(\mathbf{z}^{(n)}; \lambda)$$

where $\mathbb{H}[\pi]$ is the entropy of π , MCGD results in inclusive KL minimization as

$$\mathbb{E}_\Pi[f(\lambda, \eta)] = d_{\text{KL}}(\pi \parallel q(\cdot; \lambda)), \quad \text{and} \quad \mathbb{E}_\Pi[\mathbf{g}(\lambda, \eta)] = \nabla_\lambda d_{\text{KL}}(\pi \parallel q(\cdot; \lambda)).$$

Proof of Proposition 1. For notational convenience, we define the shorthand

$$\pi(\mathbf{z}^{(1:N)}) = \pi(\mathbf{z}^{(1)}) \pi(\mathbf{z}^{(2)}) \times \dots \times \pi(\mathbf{z}^{(N)}).$$

Then,

$$\begin{aligned} & \mathbb{E}_\Pi[f(\lambda, \eta)] \\ &= \int \left(-\frac{1}{N} \sum_{n=1}^N \log q(\mathbf{z}^{(n)}; \lambda) - \mathbb{H}[\pi] \right) \pi(\mathbf{z}^{(1:N)}) d\mathbf{z}^{(1:N)} \\ &= \int \left(-\frac{1}{N} \sum_{n=1}^N \log q(\mathbf{z}^{(n)}; \lambda) \right) \pi(\mathbf{z}^{(1:N)}) d\mathbf{z}^{(1:N)} - \mathbb{H}[\pi] \quad \text{Pulled out constant} \\ &= \frac{1}{N} \sum_{n=1}^N \left\{ \int (-\log q(\mathbf{z}^{(n)}; \lambda)) \pi(\mathbf{z}^{(1:N)}) d\mathbf{z}^{(1:N)} \right\} - \mathbb{H}[\pi] \quad \text{Swapped integral and sum} \\ &= \frac{1}{N} \sum_{n=1}^N \int (-\log q(\mathbf{z}^{(n)}; \lambda)) \pi(\mathbf{z}^{(n)}) d\mathbf{z}^{(n)} - \mathbb{H}[\pi] \quad \text{Marginalized } \mathbf{z}^{(m)} \text{ for all } m \neq n \\ &= \frac{1}{N} \sum_{n=1}^N \int (-\log q(\mathbf{z}^{(n)}; \lambda) + \log \pi(\mathbf{z}^{(n)})) \pi(\mathbf{z}^{(n)}) d\mathbf{z}^{(n)} \quad \text{Definition of entropy } \mathbb{H}[\pi] \\ &= \frac{1}{N} \sum_{n=1}^N \int \pi(\mathbf{z}^{(n)}) \log \frac{\pi(\mathbf{z}^{(n)})}{q(\mathbf{z}^{(n)}; \lambda)} d\mathbf{z}^{(n)} \\ &= \frac{1}{N} \sum_{n=1}^N d_{\text{KL}}(\pi \parallel q(\cdot; \lambda)) \quad \text{Definition of } d_{\text{KL}} \\ &= d_{\text{KL}}(\pi \parallel q(\cdot; \lambda)) \end{aligned} \tag{4}$$

For $\mathbb{E}_\Pi[\mathbf{g}(\lambda, \eta)]$, note that

$$\nabla_\lambda f(\lambda, \eta) = -\frac{1}{N} \sum_{n=1}^N \nabla_\lambda \log q(\mathbf{z}^{(n)}; \lambda) = -\frac{1}{N} \sum_{n=1}^N \mathbf{s}(\mathbf{z}^{(n)}; \lambda) = \mathbf{g}(\lambda, \eta). \tag{5}$$

Therefore, it suffices to show that

$$\begin{aligned} \nabla_\lambda d_{\text{KL}}(\pi \parallel q(\cdot; \lambda)) &= \nabla_\lambda \mathbb{E}_\Pi[f(\lambda, \eta)] && \text{Equation (4)} \\ &= \mathbb{E}_\Pi[\nabla_\lambda f(\lambda, \eta)] && \text{Leibniz derivative rule} \\ &= \mathbb{E}_\Pi[\mathbf{g}(\lambda, \eta)]. && \text{Equation (5)} \end{aligned}$$

□

Proposition 2. The maximum importance weight $w^* = \sup_{\mathbf{z}} w(\mathbf{z}) = \sup_{\mathbf{z}} \pi(\mathbf{z}) / q(\mathbf{z}; \lambda)$ is bounded below exponentially by the KL divergence as

$$\exp(d_{\text{KL}}(\pi \parallel q(\cdot; \lambda))) < w^*.$$

Proof of Proposition 2.

$$\begin{aligned}
d_{\text{KL}}(\pi \parallel q(\cdot; \lambda)) &= \mathbb{E}_{\mathbf{z} \sim \pi(\cdot)} \left[\log \frac{\pi(\mathbf{z})}{q(\mathbf{z}; \lambda)} \right] \quad \text{Definition of } d_{\text{KL}} \\
&\leq \log \mathbb{E}_{\mathbf{z} \sim \pi(\cdot)} \left[\frac{\pi(\mathbf{z})}{q(\mathbf{z}; \lambda)} \right] \quad \text{Jensen's inequality} \\
&\leq \log \mathbb{E}_{\mathbf{z} \sim \pi(\cdot)} [w^*] \\
&= \log w^*.
\end{aligned}$$

□

Lemma 1. For the probability measures p_1, \dots, p_N and q_1, \dots, q_N defined on a measurable space (X, \mathcal{A}) and an arbitrary set $A \in \mathcal{A}$,

$$\begin{aligned}
&\left| \int_{A^N} p_1(dx_1) p_2(dx_2) \times \dots \times p_N(dx_N) - q_1(dx_1) q_2(dx_2) \times \dots \times q_N(dx_N) \right| \\
&\leq \sum_{n=1}^N \left| \int_A p_n(dx_n) - q_n(dx_n) \right|
\end{aligned}$$

Proof of Lemma 1. By using the following shorthand notations

$$\begin{aligned}
p_{(1:N)}(dx_{(1:N)}) &= p_1(dx_1) p_2(dx_2) \times \dots \times p_N(dx_N) \\
q_{(1:N)}(dx_{(1:N)}) &= q_1(dx_1) q_2(dx_2) \times \dots \times q_N(dx_N),
\end{aligned}$$

the result follows from induction as

$$\begin{aligned}
&\left| \int_{A^N} p_{(1:N)}(dx_{(1:N)}) - q_{(1:N)}(dx_{(1:N)}) \right| \\
&= \left| \left(\int_A p_1(dx_1) - q_1(dx_1) \right) \int_{A^{N-1}} p_{(2:N)}(dx_{(2:N)}) \right. \\
&\quad \left. + \int_A q_1(dx_1) \left(\int_{A^{N-1}} p_{(2:N)}(dx_{(2:N)}) - q_{(2:N)}(dx_{(2:N)}) \right) \right| \\
&\leq \left| \int_A p_1(dx_1) - q_1(dx_1) \right| \int_{A^{N-1}} p_{(2:N)}(dx_{(2:N)}) \\
&\quad + \int_A q_1(dx_1) \left| \int_{A^{N-1}} p_{(2:N)}(dx_{(2:N)}) - q_{(2:N)}(dx_{(2:N)}) \right| \quad \text{Triangle inequality} \\
&\leq \left| \int_A p_1(dx_1) - q_1(dx_1) \right| \\
&\quad + \left| \int_{A^{N-1}} p_{(2:N)}(dx_{(2:N)}) - q_{(2:N)}(dx_{(2:N)}) \right|. \quad \text{Applied } p_n(A), q_n(A) \leq 1
\end{aligned}$$

□

Theorem 1. MSC (Naesseth *et al.*, 2020) is obtained by defining

$$P_\lambda^n(\eta, d\eta') = K_\lambda^n(\mathbf{z}, d\mathbf{z}')$$

with $\eta_t = \mathbf{z}_t$ where $K_\lambda(\mathbf{z}, \cdot)$ is the CIS kernel with $q_{\text{def.}}(\cdot; \lambda)$ as its proposal distribution. Then, given Assumption 2 and 3, the mixing rate and the gradient bounds are given as

$$d_{\text{TV}}(P_\lambda^n(\eta, \cdot), \Pi) \leq \left(1 - \frac{N-1}{2w^* + N-2} \right)^n \quad \text{and} \quad \mathbb{E} [\|\mathbf{g}(\lambda, \eta)\|_*^2 | \mathcal{F}_t] \leq L^2,$$

where $w^* = \sup_{\mathbf{z}} \pi(\mathbf{z}) / q_{\text{def.}}(\mathbf{z}; \lambda)$.

Proof of Theorem 1. MSC is described in Algorithm 3. At each iteration, it performs a single MCMC transition with the CIS kernel where it internally uses N proposals.

Ergodicity of the Markov Chain The ergodic convergence rate of P_λ is equal to that of K_λ , the CIS kernel proposed by Naesseth *et al.* (2020). Although not mentioned by Naesseth *et al.* (2020), this kernel has been previously proposed as the iterated sequential importance resampling (i-SIR) by Andrieu *et al.* (2018) with its corresponding geometric convergence rate as

$$d_{\text{TV}}(P_\lambda^n(\eta, \cdot), \Pi) = d_{\text{TV}}(K_\lambda^n(\mathbf{z}, \cdot), \pi) \leq \left(1 - \frac{N-1}{2w^* + N-2}\right)^n.$$

Bound on the Gradient Variance The bound on the gradient variance is straightforward given Assumption 3. For simplicity, we denote the rejection state as $\mathbf{z}^{(1)} = \mathbf{z}_{t-1}$. Then,

$$\begin{aligned} & \mathbb{E} [\|\mathbf{g}(\lambda, \eta)\|_*^2 | \mathcal{F}_t] \\ &= \mathbb{E} [\|\mathbf{g}(\lambda, \eta)\|_*^2 | \mathcal{F}_t] \\ &= \mathbb{E}_{\mathbf{z} \sim K_{\lambda_{t-1}}(\mathbf{z}_{t-1}, \cdot)} [\|\mathbf{s}(\lambda; \mathbf{z})\|_*^2 | \lambda_{t-1}, \mathbf{z}_{t-1}] \\ &= \int \sum_{n=1}^N \frac{w(\mathbf{z}^{(n)})}{\sum_{m=1}^N w(\mathbf{z}^{(m)})} \|\mathbf{s}(\cdot; \mathbf{z}^{(n)})\|_*^2 \prod_{n=2}^N q(d\mathbf{z}^{(n)}; \lambda_{t-1}) \quad (\text{Andrieu et al., 2018}) \\ &\leq L^2 \int \sum_{n=1}^N \frac{w(\mathbf{z}^{(n)})}{\sum_{m=1}^N w(\mathbf{z}^{(m)})} \prod_{n=2}^N q(d\mathbf{z}^{(n)}; \lambda_{t-1}) \quad \text{Assumption 3} \\ &= L^2 \int \prod_{n=2}^N q(d\mathbf{z}^{(n)}; \lambda_{t-1}) \quad \text{The sum of weights is 1} \\ &= L^2. \end{aligned}$$

□

Lemma 2. For $w^* = \sup_{\mathbf{z}} w(\mathbf{z}), \lambda(\cdot)$ in Equation (3) is bounded as

$$\max\left(1 - \frac{1}{w}, 0\right) \leq \lambda(w) \leq 1 - \frac{1}{w^*}.$$

Proof of Lemma 2. The proof can be found in the proof of Theorem 3 of Smith & Tierney (1996). □

Lemma 3. For $w^* = \sup_{\mathbf{z}} w(\mathbf{z}), T_n(\cdot)$ in Equation (3) is bounded as

$$T_n(w) \leq \frac{n}{w} \left(1 - \frac{1}{w^*}\right)^{n-1}.$$

Proof of Lemma 3.

$$\begin{aligned} T_n(w) &= \int_w^\infty \frac{n}{v^2} \lambda^{n-1}(v) dv \quad \text{Equation (3)} \\ &\leq \int_w^\infty \frac{n}{v^2} \left(1 - \frac{1}{w^*}\right)^{n-1} dv \quad \text{Lemma 2} \\ &= n \left(1 - \frac{1}{w^*}\right)^{n-1} \int_w^\infty \frac{1}{v^2} dv \quad \text{Pulled out constant} \\ &= n \left(1 - \frac{1}{w^*}\right)^{n-1} \left(-\frac{1}{v}\Big|_w^\infty\right) \quad \text{Solved indefinite integral} \\ &= \frac{n}{w} \left(1 - \frac{1}{w^*}\right)^{n-1}. \end{aligned}$$

This upper bound is in general difficult to improve unless we impose stronger assumptions on π and q . □

Lemma 4. For a positive test function $f : \mathcal{Z} \rightarrow \mathbb{R}^+$, the estimate of a π -invariant independent Metropolis-Hastings kernel with a proposal q is bounded as

$$\mathbb{E}_{K^n(\mathbf{z}, \cdot)} [f | \mathbf{z}] \leq n \left(1 - \frac{1}{w^*}\right)^{n-1} \mathbb{E}_q [f] + \left(1 - \frac{1}{w^*}\right)^n f(\mathbf{z})$$

where $w(\mathbf{z}) = \pi(\mathbf{z})/q(\mathbf{z})$ and $w^* = \sup_{\mathbf{z}} w(\mathbf{z})$.

Proof of Lemma 4.

$$\begin{aligned} & \mathbb{E}_{K^n(\mathbf{z}, \cdot)} [f | \mathbf{z}] \\ &= \int T_n(w(\mathbf{z}) \vee w(\mathbf{z}')) f(\mathbf{z}') \pi(\mathbf{z}') d\mathbf{z}' + \lambda^n(w(\mathbf{z})) f(\mathbf{z}) && \text{Equation (2)} \\ &\leq \int \frac{n}{w(\mathbf{z}) \vee w(\mathbf{z}')} \left(1 - \frac{1}{w^*}\right)^{n-1} f(\mathbf{z}') \pi(\mathbf{z}') d\mathbf{z}' + \lambda^n(w(\mathbf{z})) f(\mathbf{z}) && \text{Lemma 3} \\ &\leq \int \frac{n}{w(\mathbf{z}')} \left(1 - \frac{1}{w^*}\right)^{n-1} f(\mathbf{z}') \pi(\mathbf{z}') d\mathbf{z}' + \lambda^n(w(\mathbf{z})) f(\mathbf{z}) && \frac{1}{w(\mathbf{z}) \vee w(\mathbf{z}')} \leq \frac{1}{w(\mathbf{z}')} \\ &= n \left(1 - \frac{1}{w^*}\right)^{n-1} \int \frac{1}{w(\mathbf{z}')} f(\mathbf{z}') \pi(\mathbf{z}') d\mathbf{z}' + \lambda^n(w(\mathbf{z})) f(\mathbf{z}) && \text{Pulled out constant} \\ &= n \left(1 - \frac{1}{w^*}\right)^{n-1} \int f(\mathbf{z}') q(\mathbf{z}') d\mathbf{z}' + \lambda^n(w(\mathbf{z})) f(\mathbf{z}) && \text{Definition of } w(\mathbf{z}) \\ &\leq n \left(1 - \frac{1}{w^*}\right)^{n-1} \int f(\mathbf{z}') q(\mathbf{z}') d\mathbf{z}' + \left(1 - \frac{1}{w^*}\right)^n f(\mathbf{z}) && \text{Lemma 2} \\ &= n \left(1 - \frac{1}{w^*}\right)^{n-1} \mathbb{E}_q [f] + \left(1 - \frac{1}{w^*}\right)^n f(\mathbf{z}). \end{aligned}$$

□

Theorem 2. JSA (Ou & Song, 2020) is obtained by defining

$$P_\lambda^n(\eta, d\eta') = K_\lambda^{N(n-1)+1}(\mathbf{z}^{(1)}, d\mathbf{z}'^{(1)}) K_\lambda^{N(n-1)+2}(\mathbf{z}^{(2)}, d\mathbf{z}'^{(2)}) \cdots K_\lambda^{N(n-1)+N}(\mathbf{z}^{(N)}, d\mathbf{z}'^{(N)})$$

with $\eta_t = (\mathbf{z}_t^{(1)}, \mathbf{z}_t^{(2)}, \dots, \mathbf{z}_t^{(N)})$. Then, given Assumption 2 and 3, the mixing rate and the gradient variance bounds are

$$d_{\text{TV}}(P_\lambda^n(\eta, \cdot), \Pi) \leq C(\rho, N) \rho^{nN} \quad \text{and} \quad \mathbb{E} [\|\mathbf{g}(\lambda, \eta)\|_*^2 | \mathcal{F}_t] \leq L^2 \left[\frac{1}{2} + \frac{3}{2} \frac{1}{N} + \mathcal{O}(1/w^*) \right],$$

where $w^* = \sup_{\mathbf{z}} \pi(\mathbf{z})/q_{\text{def.}}(\mathbf{z}; \lambda)$ and $C(\rho, N) > 0$ is a finite constant depending on ρ and N .

Proof of Theorem 2. JSA is described in Algorithm 4. At each iteration, it performs N MCMC transitions, and uses the N samples to estimate the gradient.

Ergodicity of the Markov Chain The state transitions of the Markov chain samples $\mathbf{z}^{(1:N)}$ are visualized as

	$\mathbf{z}_t^{(1)}$	$\mathbf{z}_t^{(2)}$	$\mathbf{z}_t^{(3)}$...	$\mathbf{z}_t^{(N)}$
$t = 1$	$K_{\lambda_1}(\mathbf{z}_0, d\mathbf{z}_1^{(1)})$	$K_{\lambda_1}^2(\mathbf{z}_0, d\mathbf{z}_1^{(2)})$	$K_{\lambda_1}^3(\mathbf{z}_0, d\mathbf{z}_1^{(3)})$...	$K_{\lambda_1}^N(\mathbf{z}_0, d\mathbf{z}_1^{(N)})$
$t = 2$	$K_{\lambda_2}^{N+1}(\mathbf{z}_0, d\mathbf{z}_2^{(1)})$	$K_{\lambda_2}^{N+2}(\mathbf{z}_0, d\mathbf{z}_2^{(2)})$	$K_{\lambda_2}^{N+3}(\mathbf{z}_0, d\mathbf{z}_2^{(3)})$...	$K_{\lambda_2}^{2N}(\mathbf{z}_0, d\mathbf{z}_2^{(N)})$
\vdots					\vdots
$t = k$	$K_{\lambda_k}^{(k-1)N+1}(\mathbf{z}_0, d\mathbf{z}_k^{(1)})$	$K_{\lambda_k}^{(k-1)N+2}(\mathbf{z}_0, d\mathbf{z}_k^{(2)})$	$K_{\lambda_k}^{(k-1)N+3}(\mathbf{z}_0, d\mathbf{z}_k^{(3)})$...	$K_{\lambda_k}^{(k-1)N+N}(\mathbf{z}_0, d\mathbf{z}_k^{(N)})$

where $K_\lambda(\mathbf{z}, \cdot)$ is an IMH kernel. Therefore, the n -step transition kernel for the vector of the Markov-chain samples $\eta = \mathbf{z}^{(1:N)}$ is represented as

$$P_\lambda^n(\eta, d\eta') = K_\lambda^{N(n-1)+1}(\mathbf{z}_1, d\mathbf{z}'_1) K_\lambda^{N(n-1)+2}(\mathbf{z}_2, d\mathbf{z}'_2) \cdots K_\lambda^{N(n-1)+N}(\mathbf{z}_N, d\mathbf{z}'_N).$$

Now, the convergence in total variation $d_{\text{TV}}(\cdot, \cdot)$ can be shown to decrease geometrically as

$$\begin{aligned}
& d_{\text{TV}}(P_\lambda^n(\eta, \cdot), \Pi) \\
&= \sup_A |\Pi(A) - P^n(\eta, A)| && \text{Definition of } d_{\text{TV}} \\
&\leq \sup_A \left| \int_A \pi(dz'^{(1)}) \times \dots \times \pi(dz'^{(N)}) \right. \\
&\quad \left. - K_\lambda^{(n-1)N+1}(z^{(1)}, dz'^{(1)}) \times \dots \times K_\lambda^{nN}(z^{(N)}, dz'^{(N)}) \right| \\
&\leq \sup_A \sum_{n=1}^N \left| \int_A \pi(dz^{(n)}) - K_\lambda^{(n-1)N+n}(z^{(n)}, dz'^{(n)}) \right| && \text{Lemma 1} \\
&= \sum_{n=1}^N d_{\text{TV}}(K_\lambda^{(n-1)N+n}(z^{(n)}, \cdot), \pi) && \text{Definition of } d_{\text{TV}} \\
&\leq \sum_{n=1}^N \rho^{(n-1)N+n} && \text{Geometric ergodicity} \\
&= \rho^{nN} \rho^{-N} \frac{\rho - \rho^{N+1}}{1 - \rho} && \text{Solved sum} \\
&= \frac{\rho(1 - \rho^N)}{\rho^N(1 - \rho)} (\rho^N)^n.
\end{aligned}$$

Although the constant depends on ρ and N , the kernel P is geometrically ergodic and converges N times faster than the base kernel K .

Bound on the Gradient Variance To analyze the variance of the gradient, we require detailed information about the n -step marginal transition kernel, which is unavailable in general for most MCMC kernels. Fortunately, specifically for the IMH kernel, Smith & Tierney (1996) have shown that the n -step marginal IMH kernel is given as Equation (2). From this, we show that

$$\begin{aligned}
& \mathbb{E} [\|\mathbf{g}(\lambda, \eta)\|_*^2 | \mathcal{F}_t] \\
&= \mathbb{E} [\|\mathbf{g}(\lambda, \eta)\|_*^2 | z_{t-1}^{(N)}, \lambda_{t-1}] \\
&= \mathbb{E} \left[\left\| \frac{1}{N} \sum_{n=1}^N \mathbf{s}(\lambda; z^{(n)}) \right\|_*^2 \middle| z_{t-1}^{(N)}, \lambda_{t-1} \right] \\
&\leq \mathbb{E} \left[\frac{1}{N^2} \sum_{n=1}^N \left\| \mathbf{s}(\lambda; z^{(n)}) \right\|_*^2 \middle| z_{t-1}^{(N)}, \lambda_{t-1} \right] && \text{Triangle inequality} \\
&= \frac{1}{N^2} \sum_{n=1}^N \mathbb{E}_{z^{(n)} \sim K^n(z_{t-1}, \cdot)} \left[\left\| \mathbf{s}(\lambda; z^{(n)}) \right\|_*^2 \middle| z_{t-1}^{(N)}, \lambda_{t-1} \right] && \text{Linearity of expectation} \\
&\leq \frac{1}{N^2} \sum_{n=1}^N n \left(1 - \frac{1}{w^*}\right)^{n-1} \mathbb{E}_{z^{(n)} \sim q_{\text{def.}(\cdot; \lambda)}} \left[\left\| \mathbf{s}(\lambda; z^{(n)}) \right\|_*^2 \right] \\
&\quad + \left(1 - \frac{1}{w^*}\right)^n \left\| \mathbf{s}(\lambda; z_{t-1}^{(N)}) \right\|_*^2 && \text{Lemma 4} \\
&\leq \frac{1}{N^2} \sum_{n=1}^N n \left(1 - \frac{1}{w^*}\right)^{n-1} L^2 + \left(1 - \frac{1}{w^*}\right)^n L^2 && \text{Assumption 3} \\
&= \frac{L^2}{N^2} \sum_{n=1}^N n \left(1 - \frac{1}{w^*}\right)^{n-1} + \left(1 - \frac{1}{w^*}\right)^n && \text{Moved constant forward}
\end{aligned}$$

$$\begin{aligned}
&= \frac{L^2}{N^2} \left[(w^*)^2 + w^* - \left(1 - \frac{1}{w^*}\right)^N ((w^*)^2 + w^* + N w^*) \right] \quad \text{Solved sum} \\
&= \frac{L^2}{N^2} \left[\frac{1}{2} N^2 + \frac{3}{2} N + \mathcal{O}(1/w^*) \right] \quad \text{Laurent series expansion at } w^* \rightarrow \infty \\
&= L^2 \left[\frac{1}{2} + \frac{3}{2} \frac{1}{N} + \mathcal{O}(1/w^*) \right].
\end{aligned}$$

The Laurent approximation becomes exact as $w^* \rightarrow \infty$, which is useful accurate as a consequence of Proposition 2. \square

Theorem 3. pMCSA, our proposed scheme, is obtained by setting

$$P_\lambda^n(\eta, d\eta') = K_\lambda^n(\mathbf{z}^{(1)}, d\mathbf{z}'^{(1)}) K_\lambda^n(\mathbf{z}^{(2)}, d\mathbf{z}'^{(2)}) \cdots K_\lambda^n(\mathbf{z}^{(N)}, d\mathbf{z}'^{(N)})$$

with $\eta = (\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(N)})$. Then, given Assumption 2 and 3, the mixing rate and the gradient variance bounds are

$$d_{\text{TV}}(P_\lambda^n(\eta, \cdot), \Pi) \leq C(N) \rho^n \quad \text{and} \quad \mathbb{E}[\|\mathbf{g}(\lambda, \eta)\|_*^2 | \mathcal{F}_t] \leq L^2 \left[\frac{1}{N} + \frac{1}{N} \left(1 - \frac{1}{w^*}\right) \right],$$

where $w^* = \sup_{\mathbf{z}} \pi(\mathbf{z}) / q_{\text{def.}}(\mathbf{z}; \lambda)$ and C is some positive constant depending on N .

Proof of Theorem 3. Our proposed scheme, pMCSA, is described in Algorithm 4. At each iteration, our scheme performs a single MCMC transition for each of the N samples, or chains, to estimate the gradient. Similarly to JSA, we use the IMH kernel K_λ .

Ergodicity of the Markov Chain Since our kernel operates the same MCMC kernel K_λ for each of the N parallel Markov chains, the n -step marginal kernel P_λ can be represented as

$$P_\lambda^n(\eta, d\eta') = K_\lambda^n(\mathbf{z}^{(1)}, d\mathbf{z}'^{(1)}) K_\lambda^n(\mathbf{z}^{(2)}, d\mathbf{z}'^{(2)}) \cdots K_\lambda^n(\mathbf{z}^{(N)}, d\mathbf{z}'^{(N)}).$$

Then, the convergence in total variation $d_{\text{TV}}(\cdot, \cdot)$ can be shown to decrease geometrically as

$$\begin{aligned}
&d_{\text{TV}}(K_\lambda^k(\eta, \cdot), \Pi) \\
&= \sup_A |\Pi(A) - P_\lambda^n(\eta, A)| \quad \text{Definition of } d_{\text{TV}} \\
&\leq \sup_A \left| \int_A \pi(d\mathbf{z}'_1) \cdot \dots \cdot \pi(d\mathbf{z}'_N) \right. \\
&\quad \left. - K_\lambda^n(\mathbf{z}_1, d\mathbf{z}'_1) \cdot \dots \cdot K_\lambda^n(\mathbf{z}_N, d\mathbf{z}'_N) \right| \\
&\leq \sup_A \sum_{n=1}^N \left| \int_A \pi(d\mathbf{z}'_k) - K_\lambda^n(\mathbf{z}_n, d\mathbf{z}'_n) \right| \quad \text{Lemma 1} \\
&= \sum_{n=1}^N d_{\text{TV}}(K_\lambda^n(\mathbf{z}_n, \cdot), \pi) \quad \text{Definition of TV distance} \\
&\leq \sum_{n=1}^N \rho^n \quad \text{Geometric ergodicity} \\
&= N \rho^k \quad \text{Solved sum.}
\end{aligned}$$

Bound on the Gradient Variance The bound on the gradient variance can be derived in similar manner to JSA as

$$\begin{aligned}
&\mathbb{E}[\|\mathbf{g}(\lambda, \eta)\|_*^2 | \mathcal{F}_t] \\
&= \mathbb{E}[\|\mathbf{g}(\lambda, \eta)\|_*^2 | \mathbf{z}_{t-1}^{(1:N)}, \lambda_{t-1}] \\
&= \mathbb{E}\left[\left\| \frac{1}{N} \sum_{n=1}^N \mathbf{s}(\lambda; \mathbf{z}_n) \right\|_*^2 | \mathbf{z}_{t-1}^{(1:N)}, \lambda_{t-1}\right]
\end{aligned}$$

$$\begin{aligned}
&\leq \mathbb{E} \left[\frac{1}{N^2} \sum_{n=1}^N \| \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_n) \|_*^2 \middle| \mathbf{z}_{t-1}^{(1:N)}, \boldsymbol{\lambda}_{t-1} \right] && \text{Triangle inequality} \\
&= \frac{1}{N^2} \sum_{n=1}^N \mathbb{E}_{\mathbf{z}_n \sim K_{\boldsymbol{\lambda}_{t-1}}(\mathbf{z}_{t-1}^{(n)}, \cdot)} \left[\| \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_n) \|_*^2 \middle| \mathbf{z}_{t-1}^{(1:N)}, \boldsymbol{\lambda}_{t-1} \right] && \text{Linearity of expectation} \\
&\leq \frac{1}{N^2} \sum_{n=1}^N \mathbb{E}_{\mathbf{z}_n \sim q_{\text{def.}(\cdot; \boldsymbol{\lambda}_{t-1})}} \left[\| \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_n) \|_*^2 \right] + \left(1 - \frac{1}{w^*}\right) \left\| \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_{t-1}^{(n)}) \right\|_*^2 && \text{Lemma 4} \\
&\leq \frac{1}{N^2} \sum_{n=1}^N L^2 + \left(1 - \frac{1}{w^*}\right) L^2 && \text{Assumption 3} \\
&= \frac{L^2}{N^2} \sum_{n=1}^N 1 + \left(1 - \frac{1}{w^*}\right) && \text{Moved constant forward} \\
&= L^2 \left[\frac{1}{N} + \frac{1}{N} \left(1 - \frac{1}{w^*}\right) \right]. && \text{Solved sum}
\end{aligned}$$

□

E Additional Experimental Results

E.1 Bayesian Neural Network Regression

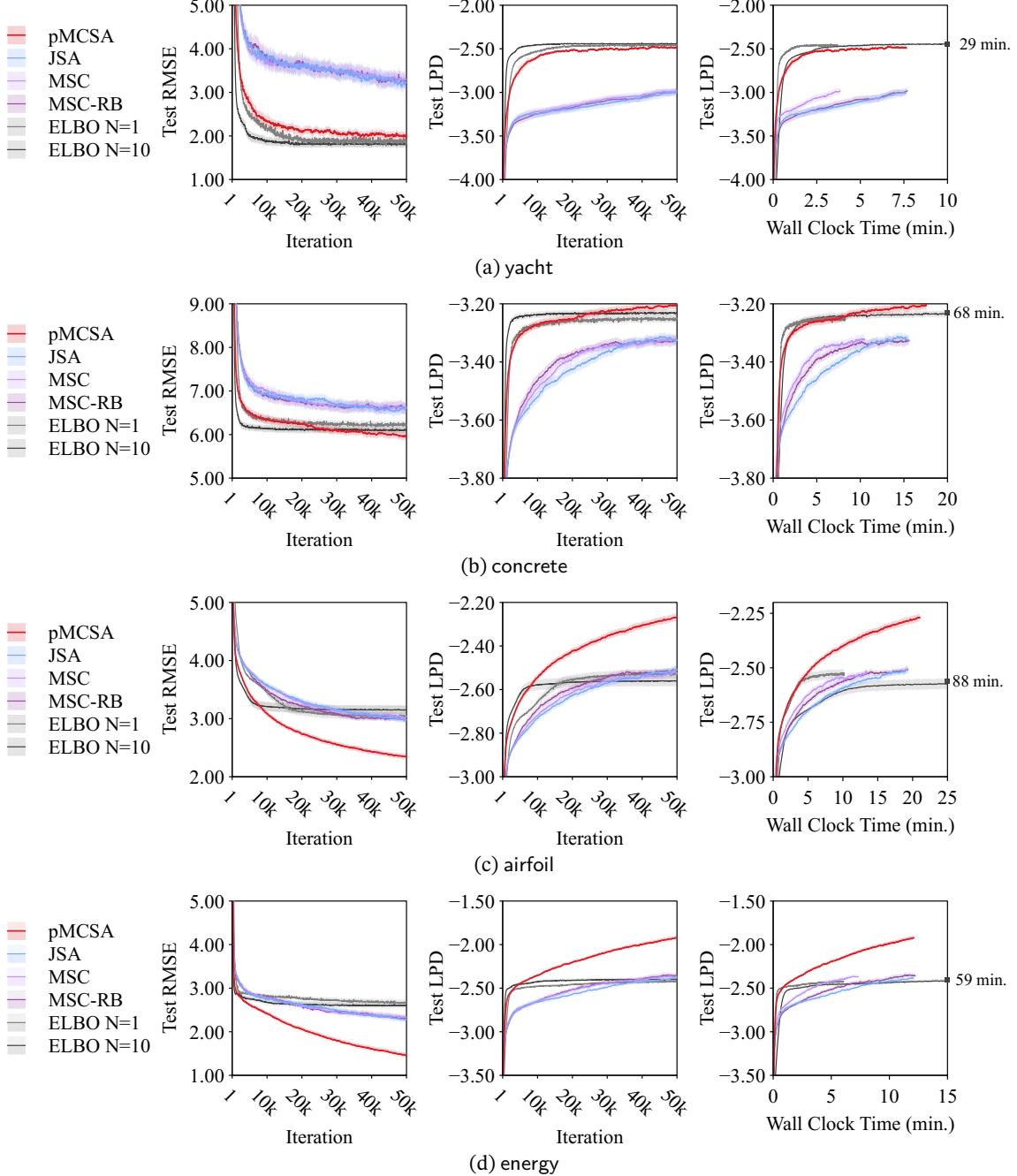


Figure 3: **Test root-mean-square error (RMSE) and test log predictive density (LPD) on Bayesian neural network regression.** The grey squares (■) mark the performance of ELBO $N = 10$ at the wall clock time shown next to it. The error bands show the 95% bootstrap confidence intervals obtained from 20 independent 90% train-test splits.

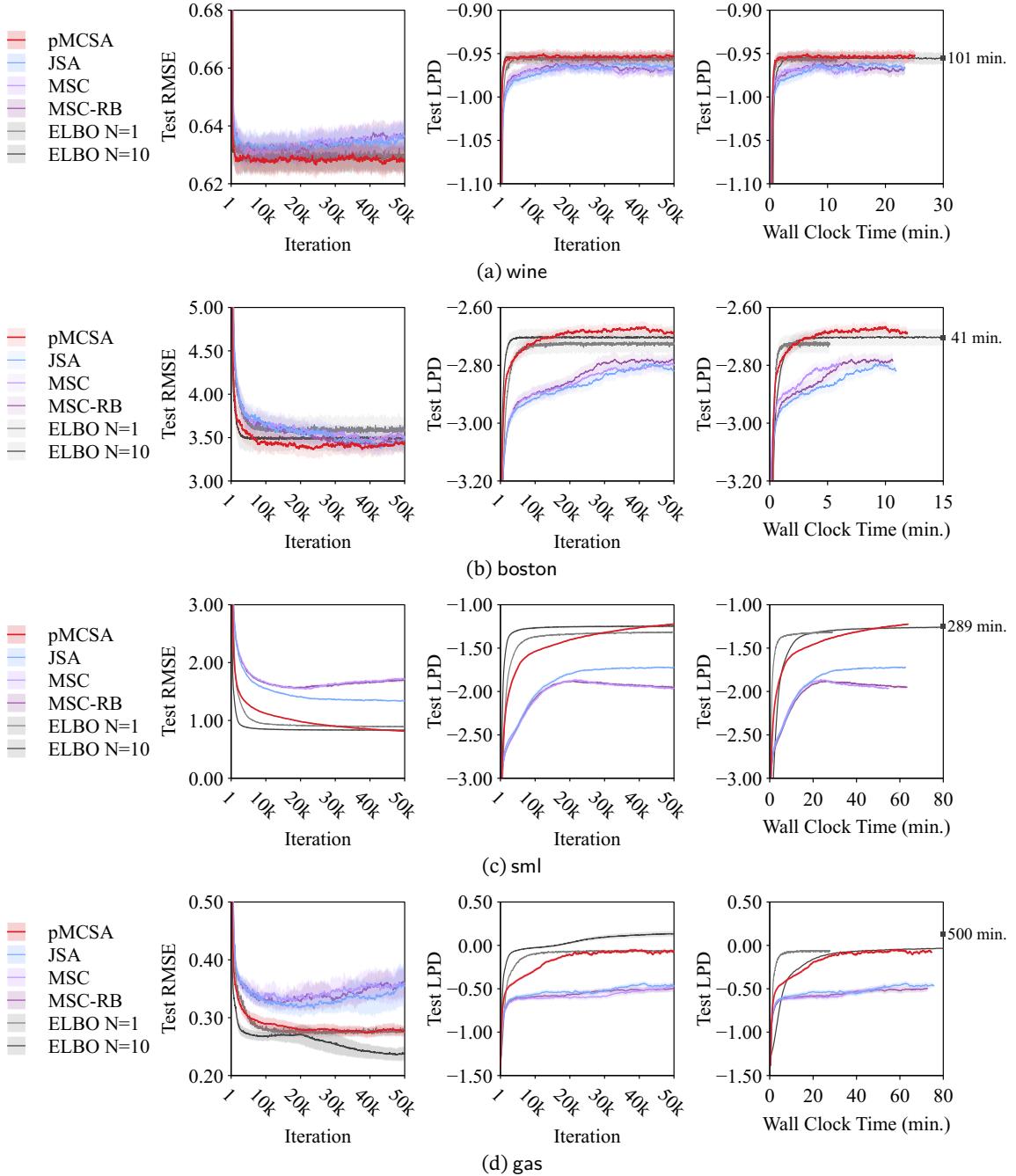


Figure 4: **(continued) Test root-mean-square error (RMSE) and test log predictive density (LPD) on Bayesian neural network regression.** The grey squares (\blacksquare) mark the performance of ELBO $N = 10$ at the wall clock time shown next to it. The error bands show the 95% bootstrap confidence intervals obtained from 20 independent 90% train-test splits.

E.2 Robust Gaussian Process Regression

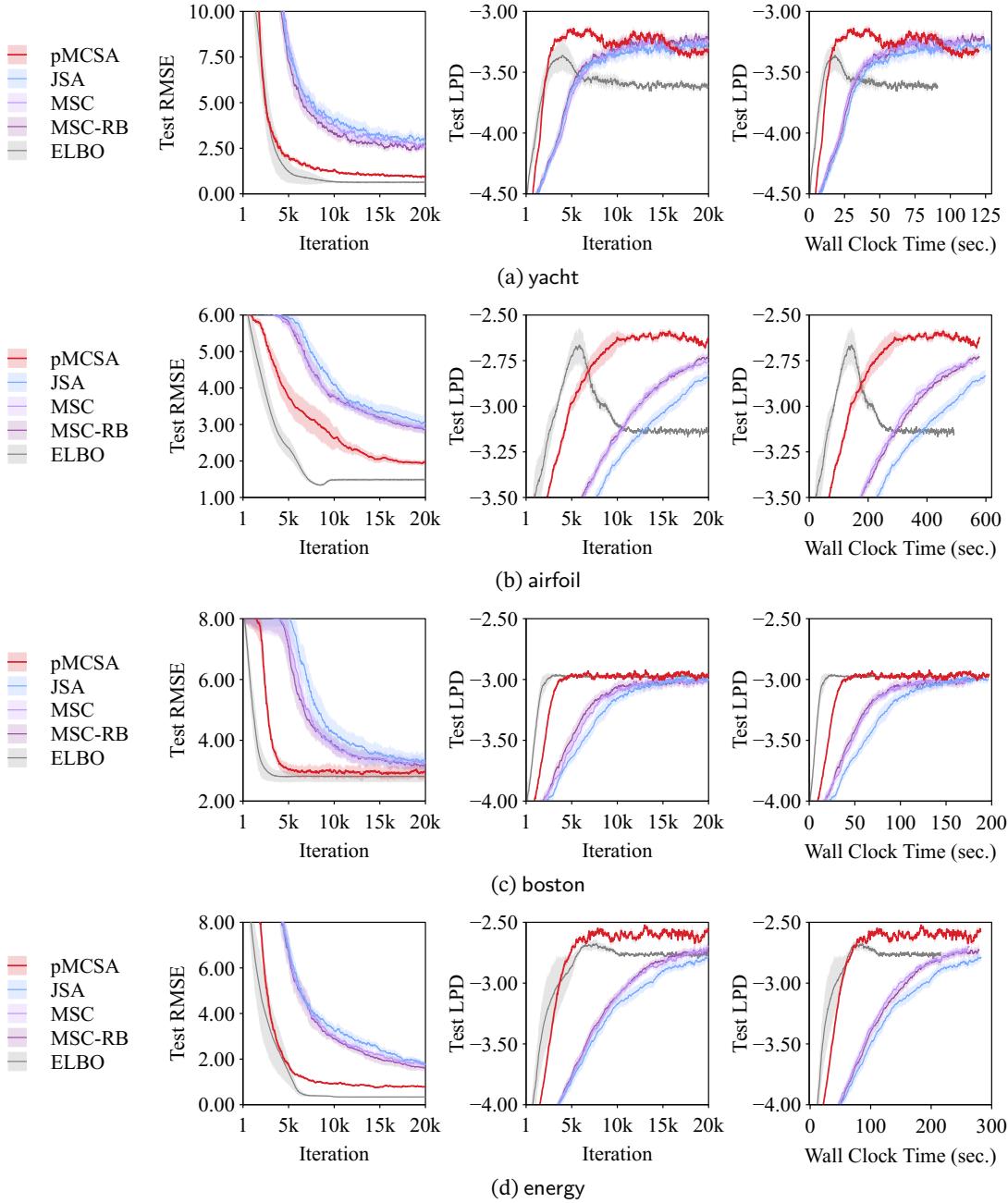


Figure 5: **Test root-mean-square error (RMSE) and test log predictive density (LPD) on robust Gaussian process regression.** The error bands shows the 95% bootstrap confidence interval obtained from 20 repetitions.

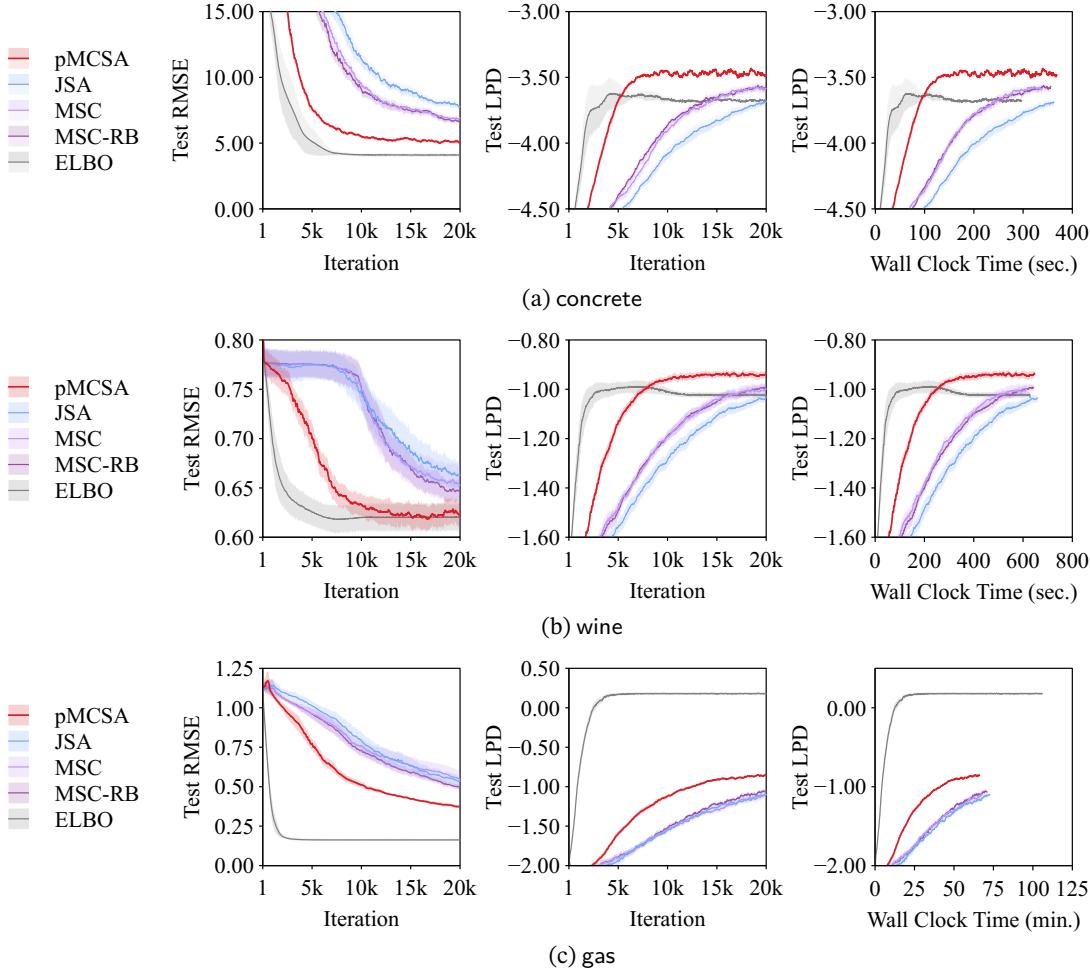


Figure 6: **(continued) Test root-mean-square error (RMSE) and test log predictive density (LPD) on robust Gaussian process regression.** The error bands show the 95% bootstrap confidence intervals obtained from 20 independent 90% train-test splits.