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

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 Monte Carlo Kernels

Algorithm 2: 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 3: Independent Metropolis-Hastings Kernel

Input: previous sample \mathbf{z}_{t-1} ,
previous parameter λ_{t-1} ,

$\mathbf{z}^* \sim q_{\text{def.}}(\mathbf{z}; \lambda_{t-1})$
 $\tilde{w}(\mathbf{z}) = p(\mathbf{z}, \mathbf{x})/q_{\text{def.}}(\mathbf{z}; \lambda_{t-1})$
 $\alpha = \min(\tilde{w}(\mathbf{z}^*)/\tilde{w}(\mathbf{z}_{t-1}), 1)$
 $u \sim \text{Uniform}(0, 1)$
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 4: Markovian Score Climbing

Input: 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 5: Joint Stochastic Approximation

Input: 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

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

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

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

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

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

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

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

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

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

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

Proof. 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,

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

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

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

Therefore, it suffices to show that

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

□

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

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

Proof.

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

□

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. 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

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

□

Lemma 2. Let \mathbf{g} be a vector-valued, biased estimator of $\boldsymbol{\mu}$, where the bias is denoted as $\text{Bias}[\mathbf{g}] = \mathbb{E}\mathbf{g} - \boldsymbol{\mu}$ and the mean-squared error is denoted as $\text{MSE}[\mathbf{g}] = \mathbb{E}\|\mathbf{g} - \boldsymbol{\mu}\|_2^2$. Then, the second moment of \mathbf{g} is bounded as

$$\begin{aligned}
\textcircled{1} \quad \mathbb{E}\|\mathbf{g}\|_2^2 &\leq \mathbb{V}\mathbf{g} + \text{Bias}[\mathbf{g}]^2 + 2 \text{Bias}[\mathbf{g}] \|\boldsymbol{\mu}\|_2 + \|\boldsymbol{\mu}\|_2^2, \\
\textcircled{2} \quad \mathbb{E}\|\mathbf{g}\|_2^2 &\leq \text{MSE}[\mathbf{g}] + 2 \text{Bias}[\mathbf{g}] \|\boldsymbol{\mu}\|_2 + \|\boldsymbol{\mu}\|_2^2,
\end{aligned}$$

where $\mathbb{V}\mathbf{g} = \mathbb{E}\|\mathbf{g} - \mathbb{E}\mathbf{g}\|_2^2$ is the variance of the estimator.

Proof. ❶ follows from the decomposition

$$\begin{aligned}
\mathbb{E} \left[\|\mathbf{g}\|_2^2 \right] &= \mathbb{V}\mathbf{g} + \|\mathbb{E}\mathbf{g}\|_2^2 \\
&= \mathbb{V}\mathbf{g} + \|\mathbb{E}\mathbf{g} - \boldsymbol{\mu} + \boldsymbol{\mu}\|_2^2 \\
&= \mathbb{V}\mathbf{g} + \|\mathbb{E}\mathbf{g} - \boldsymbol{\mu}\|_2^2 + 2(\mathbb{E}\mathbf{g} - \boldsymbol{\mu})^\top \boldsymbol{\mu} + \|\boldsymbol{\mu}\|_2^2 && \text{Expanded quadratic} \\
&\leq \mathbb{V}\mathbf{g} + \|\mathbb{E}\mathbf{g} - \boldsymbol{\mu}\|_2^2 + 2\|\mathbb{E}\mathbf{g} - \boldsymbol{\mu}\|_2 \|\boldsymbol{\mu}\|_2 + \|\boldsymbol{\mu}\|_2^2 && \text{Cauchy-Schwarz inequality} \\
&= \mathbb{V}\mathbf{g} + \text{Bias}[\mathbf{g}]^2 + 2\text{Bias}[\mathbf{g}] \|\boldsymbol{\mu}\|_2 + \|\boldsymbol{\mu}\|_2^2. && \text{Definition of bias}
\end{aligned}$$

Meanwhile, by the well-known bias-variance decomposition formula of the mean-squared error, ❷ directly follows from ❶. \square

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

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

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

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

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

Proof. MSC is described in Algorithm 4. At each iteration, it performs a single MCMC transition with the CIS kernel where it internally uses N proposals. That is,

$$\begin{aligned}
\mathbf{z}_t \mid \mathbf{z}_{t-1}, \boldsymbol{\lambda}_{t-1} &\sim K_{\boldsymbol{\lambda}_{t-1}}(\mathbf{z}_{t-1}, \cdot) \\
\mathbf{g}_{t, \text{MSC}} &= -\mathbf{s}(\boldsymbol{\lambda}, \mathbf{z}_t),
\end{aligned}$$

where $K_{\boldsymbol{\lambda}_{t-1}}$ is the CIS kernel using $q_{\text{def}}(\cdot; \boldsymbol{\lambda}_{t-1})$.

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

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

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

$$\begin{aligned}
&\mathbb{E} \left[\|\mathbf{g}_{t, \text{MSC}}\|^2 \mid \mathcal{F}_{t-1} \right] \\
&= \mathbb{E}_{\mathbf{z}_t \sim K_{\boldsymbol{\lambda}_{t-1}}(\mathbf{z}_{t-1}, \cdot)} \left[\|\mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t)\|^2 \mid \boldsymbol{\lambda}_{t-1}, \mathbf{z}_{t-1} \right] \\
&= \int \sum_{n=1}^N \frac{w(\mathbf{z}_t^{(n)})}{\sum_{m=1}^N w(\mathbf{z}_t^{(m)})} \|\mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(n)})\|^2 \prod_{n=2}^N q(d\mathbf{z}_t^{(n)}; \boldsymbol{\lambda}_{t-1}) && \text{Andrieu et al. (2018)} \\
&\leq L^2 \int \sum_{n=1}^N \frac{w(\mathbf{z}_t^{(n)})}{\sum_{m=1}^N w(\mathbf{z}_t^{(m)})} \prod_{n=2}^N q(d\mathbf{z}_t^{(n)}; \boldsymbol{\lambda}_{t-1}) && \text{Assumption 3} \\
&= L^2 \int \prod_{n=2}^N q(d\mathbf{z}_t^{(n)}; \boldsymbol{\lambda}_{t-1}) \\
&= L^2.
\end{aligned}$$

\square

Lemma 3. Let the importance weight be defined as $w(\mathbf{z}) = \pi(\mathbf{z})/q(\mathbf{z})$. The variance of the importance weights is related to the χ^2 divergence as

$$\mathbb{V}_q w(\mathbf{z}) = d_{\chi^2}(\pi \parallel q).$$

Proof.

$$\mathbb{V}_q w(\mathbf{z}) = \mathbb{E}_q \left[(w(\mathbf{z}) - \mathbb{E}_q[w(\mathbf{z})])^2 \right] = \mathbb{E}_q \left[(w(\mathbf{z}) - 1)^2 \right] = \int \left(\frac{\pi(\mathbf{z})}{q(\mathbf{z})} - 1 \right)^2 q(d\mathbf{z}) = d_{\chi^2}(\pi \parallel q).$$

□

Theorem 2. (Cardoso *et al.*, 2022) The gradient variance of MSC-RB is bounded as

$$\mathbb{E} \left[\|\mathbf{g}_{t,\text{MSC-RB}}\|_2^2 \mid \mathcal{F}_{t-1} \right] \leq 4L^2 \left[\frac{1}{N-1} d_{\chi^2}(\pi \parallel q(\cdot; \boldsymbol{\lambda}_{t-1})) + \mathcal{O}(N^{-3/2} + \gamma^{t-1}/N-1) \right] + \|\boldsymbol{\mu}\|_2^2,$$

where $\boldsymbol{\mu} = \mathbb{E}_\pi \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z})$, $\gamma = 2w^*/(2w^* + N - 2)$ is the mixing rate of the Rao-Blackwellized CIS kernel, and $d_{\chi^2}(\pi \parallel q) = \int (\pi/q - 1)^2 q(d\mathbf{z})$ is the χ^2 divergence.

Proof. Rao-Blackwellization of the CIS kernel is to reuse the importance weights $w(\mathbf{z}) = \pi(\mathbf{z})/q_{\text{def}}(\mathbf{z})$ internally used by the kernel when forming the estimator. That is, the gradient is estimated as

$$\begin{aligned} \mathbf{z}_t^{(n)} \mid \boldsymbol{\lambda}_{t-1} &\sim q(\cdot; \boldsymbol{\lambda}_{t-1}) \\ \mathbf{g}_{t,\text{MSC-RB}} &= - \sum_{n=2}^N \frac{w(\mathbf{z}_t^{(n)})}{\sum_{m=2}^N w(\mathbf{z}_t^{(m)}) + w(\mathbf{z}_{t-1})} \mathbf{s}(\mathbf{z}_t^{(n)}) + \frac{w(\mathbf{z}_{t-1})}{\sum_{m=2}^N w(\mathbf{z}_t^{(m)}) + w(\mathbf{z}_{t-1})} \mathbf{s}(\mathbf{z}_{t-1}). \end{aligned}$$

By Lemma 2, the second moment of the gradient is bounded as

$$\begin{aligned} \mathbb{E} \left[\|\mathbf{g}_{t,\text{MSC-RB}}\|_2^2 \mid \mathcal{F}_{t-1} \right] &= \text{MSE} \left[\mathbf{g}_{t,\text{MSC-RB}} \mid \mathcal{F}_{t-1} \right] + 2 \text{Bias} \left[\mathbf{g}_{t,\text{MSC-RB}} \mid \mathcal{F}_{t-1} \right]^\top \boldsymbol{\mu} + \|\boldsymbol{\mu}\|_2^2 \\ &\leq \text{MSE} \left[\mathbf{g}_{t,\text{MSC-RB}} \mid \mathcal{F}_{t-1} \right] + 2L \text{Bias} \left[\mathbf{g}_{t,\text{MSC-RB}} \mid \mathcal{F}_{t-1} \right] + \|\boldsymbol{\mu}\|_2^2. \end{aligned} \quad (6)$$

Cardoso *et al.* (2022, Theorem 3) show that the mean-squared error of this estimator, which they call bias reduced self-normalized importance sampling, is bounded as

$$\begin{aligned} \text{MSE} \left[\mathbf{g}_{t,\text{MSC-RB}} \mid \mathcal{F}_{t-1} \right] &\leq 4L^2 \left[(1 + \epsilon^2) \frac{1}{N-1} \mathbb{V}_{\mathbf{z} \sim q_{\text{def}}(\cdot; \boldsymbol{\lambda}_{t-1})} [w(\mathbf{z}) \mid \boldsymbol{\lambda}_{t-1}] \right. \\ &\quad \left. + (1 + \epsilon^{-2}) \frac{1}{N^2} (1 + w^*)^2 \right], \end{aligned}$$

for some arbitrary constant ϵ^2 . The first term is identical to the variance of an $N - 1$ -sample self-normalized importance sampling estimator (Agapiou *et al.*, 2017), while the second term is the added variance due to “rejections.”

Since the variance of the importance weights is well known to be related to the χ^2 divergence,

$$\begin{aligned} \text{MSE} \left[\mathbf{g}_{t,\text{MSC-RB}} \mid \mathcal{F}_{t-1} \right] &\leq 4L^2 \left[(1 + \epsilon^2) \frac{1}{N-1} d_{\chi^2}(\pi \parallel q(\cdot; \boldsymbol{\lambda}_{t-1})) \right. \\ &\quad \left. + (1 + \epsilon^{-2}) \frac{1}{N^2} (1 + w^*)^2 \right]. \end{aligned} \quad \text{Lemma 3}$$

For ϵ^2 , Cardoso *et al.* choose $\epsilon^2 = (N - 1)^{-1/2}$, which results in their stated bound

$$\text{MSE} \left[\mathbf{g}_{t,\text{MSC-RB}} \mid \mathcal{F}_{t-1} \right] \leq 4L^2 \left[\frac{1}{N-1} d_{\chi^2}(\pi \parallel q(\cdot; \boldsymbol{\lambda}_{t-1})) + \mathcal{O}(N^{-3/2}) \right].$$

Furthermore, they show that the bias term is bounded as

$$\text{Bias} \left[\mathbf{g}_{t,\text{MSC-RB}} \right] \leq \frac{4L}{N-1} (d_{\chi^2}(\pi \parallel q(\cdot; \boldsymbol{\lambda}_{t-1})) + 1 + w^*) \left(\frac{2w^*}{2w^* + N - 2} \right)^{t-1}.$$

Combining both the bias and the mean-squared error to Equation (6), we obtain the bound

$$\begin{aligned}
& \mathbb{E} \left[\|\mathbf{g}_{t,\text{MSC-RB}}\|_2^2 \mid \mathcal{F}_{t-1} \right] \\
& \leq 4L^2 \left[\frac{1}{N-1} d_{\chi^2}(\pi \parallel q(\cdot; \boldsymbol{\lambda}_{t-1})) \right. \\
& \quad \left. + \frac{1}{N-1} (d_{\chi^2}(\pi \parallel q(\cdot; \boldsymbol{\lambda}_{t-1})) + 1 + w^*) \left(\frac{2w^*}{2w^* + N - 2} \right)^{t-1} + \mathcal{O}(N^{-3/2}) \right] \\
& = 4L^2 \left[\frac{1 + \gamma^{t-1}}{N-1} d_{\chi^2}(\pi \parallel q(\cdot; \boldsymbol{\lambda}_{t-1})) + \frac{\gamma^{t-1}}{N-1} + \frac{\gamma^{t-1} w^*}{N-1} + \mathcal{O}(N^{-3/2}) \right].
\end{aligned}$$

□

Lemma 4. 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. The proof can be found in the proof of Theorem 3 of Smith & Tierney (1996). □

Lemma 5. 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.

$$\begin{aligned}
T_n(w) &= \int_w^\infty \frac{n}{v^2} \lambda^{n-1}(v) dv && \text{Equation (3)} \\
&\leq \int_w^\infty \frac{n}{v^2} \left(1 - \frac{1}{w^*}\right)^{n-1} dv && \text{Lemma 4} \\
&= n \left(1 - \frac{1}{w^*}\right)^{n-1} \int_w^\infty \frac{1}{v^2} dv \\
&= n \left(1 - \frac{1}{w^*}\right)^{n-1} \left(-\frac{1}{v}\Big|_w^\infty\right) \\
&= \frac{n}{w} \left(1 - \frac{1}{w^*}\right)^{n-1}.
\end{aligned}$$

□

Lemma 6. For a positive test function $f : \mathcal{Z} \rightarrow \mathbb{R}^+$, the estimate of a π -invariant independent Metropolis-Hastings kernel K with a proposal q is bounded as

$$\mathbb{E}_{K^n(\mathbf{z}, \cdot)} [f \mid \mathbf{z}] \leq n r^{n-1} \mathbb{E}_q [f] + r^n f(\mathbf{z}),$$

where $w(\mathbf{z}) = \pi(\mathbf{z})/q(\mathbf{z})$ and $r = 1 - 1/w^*$ for $w^* = \sup_{\mathbf{z}} w(\mathbf{z})$.

Proof.

$$\begin{aligned}
& \mathbb{E}_{K^n(\mathbf{z}, \cdot)} [f \mid \mathbf{z}] \\
&= \int T_n(w(\mathbf{z}) \vee w(\mathbf{z}')) f(\mathbf{z}') \pi(\mathbf{z}') d\mathbf{z}' + \lambda^n(w(\mathbf{z})) f(\mathbf{z}) && \text{Equation (2)} \\
&\leq \int \frac{n}{w(\mathbf{z}) \vee w(\mathbf{z}')} \left(1 - \frac{1}{w^*}\right)^{n-1} f(\mathbf{z}') \pi(\mathbf{z}') d\mathbf{z}' + \lambda^n(w(\mathbf{z})) f(\mathbf{z}) && \text{Lemma 5} \\
&\leq \int \frac{n}{w(\mathbf{z}')} \left(1 - \frac{1}{w^*}\right)^{n-1} f(\mathbf{z}') \pi(\mathbf{z}') d\mathbf{z}' + \lambda^n(w(\mathbf{z})) f(\mathbf{z}) && \frac{1}{w(\mathbf{z}) \vee w(\mathbf{z}')} \leq \frac{1}{w(\mathbf{z}')}
\end{aligned}$$

$$\begin{aligned}
&= n \left(1 - \frac{1}{w^*}\right)^{n-1} \int \frac{1}{w(\mathbf{z}')} f(\mathbf{z}') \pi(\mathbf{z}') d\mathbf{z}' + \lambda^n (w(\mathbf{z})) f(\mathbf{z}) \\
&= n \left(1 - \frac{1}{w^*}\right)^{n-1} \int f(\mathbf{z}') q(\mathbf{z}') d\mathbf{z}' + \lambda^n (w(\mathbf{z})) f(\mathbf{z}) && \text{Definition of } w(\mathbf{z}) \\
&\leq n \left(1 - \frac{1}{w^*}\right)^{n-1} \int f(\mathbf{z}') q(\mathbf{z}') d\mathbf{z}' + \left(1 - \frac{1}{w^*}\right)^n f(\mathbf{z}) && \text{Lemma 4} \\
&= n \left(1 - \frac{1}{w^*}\right)^{n-1} \mathbb{E}_q[f] + \left(1 - \frac{1}{w^*}\right)^n f(\mathbf{z}).
\end{aligned}$$

□

Lemma 7. Let a Π -invariant Markov chain kernel P be geometrically ergodic as

$$d_{\text{TV}}(P^n(\eta_0, \cdot), \Pi) \leq C \rho^n.$$

Furthermore, let $\hat{\mathbf{g}} = \mathbf{g}(\boldsymbol{\eta})$ with $\boldsymbol{\eta} \sim P^n(\eta_0, \cdot)$ be the estimator of $\mathbb{E}_{\Pi} \mathbf{g}$ for some function $\mathbf{g} : \mathbb{H} \rightarrow \mathbb{R}^D$ bounded as $\|\mathbf{g}\|_2 \leq L$. The bias of $\hat{\mathbf{g}}$, defined as $\text{Bias}[\hat{\mathbf{g}}] = \mathbb{E} \|\hat{\mathbf{g}} - \mathbb{E}_{\pi} \mathbf{g}\|$, is bounded as

$$\text{Bias}[\hat{\mathbf{g}}] \leq 2\sqrt{D} L C \rho^n.$$

Proof.

$$\begin{aligned}
\text{Bias}[\hat{\mathbf{g}}] &= \left\| \mathbb{E}_{P^n(\eta_0, \cdot)} \mathbf{g} - \mathbb{E}_{\Pi} \mathbf{g} \right\|_2 \\
&\leq \sqrt{D} \left\| \mathbb{E}_{P^n(\eta_0, \cdot)} \mathbf{g} - \mathbb{E}_{\Pi} \mathbf{g} \right\|_{\infty} && \text{for } \mathbf{x} \in \mathbb{R}^D, \|\mathbf{x}\|_2 \leq \sqrt{D} \|\mathbf{x}\|_{\infty} \\
&\leq \sqrt{D} L \sup_{|h| \leq 1} \left| \mathbb{E}_{P^n(\eta_0, \cdot)} h - \mathbb{E}_{\Pi} h \right| && \|\mathbf{g}\|_{\infty} \leq \|\mathbf{g}\|_2 \leq L \\
&= 2\sqrt{D} L d_{\text{TV}}(P^n(\eta_0, \cdot), \Pi) && \text{Definition of } d_{\text{TV}} \\
&= 2\sqrt{D} L C \rho^n. && \text{Geometric ergodicity}
\end{aligned}$$

□

Theorem 3. JSA (Ou & Song, 2020) is obtained by defining

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

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

$$d_{\text{TV}}(P_{\lambda}^k(\boldsymbol{\eta}, \cdot), \Pi) \leq C(r, N) r^{kN} \text{ and } \mathbb{E} \left[\|\mathbf{g}_{t, \text{JSA}}\|_2^2 \mid \mathcal{F}_{t-1} \right] \leq L^2 \left[\frac{1}{2} + \frac{3}{2} \frac{1}{N} + \mathcal{O}(1/w^* + r^t N) \right] + C_{\text{cov}} + \|\boldsymbol{\mu}\|_2^2,$$

where $\boldsymbol{\mu} = \mathbb{E}_{\pi} \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z})$, $C_{\text{cov}} = \frac{2}{N^2} \sum_{n=2}^N \sum_{m=1}^{n-1} \text{Cov}(\mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(n)}), \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(m)}) \mid \mathcal{F}_{t-1})$ is the sum of the covariance between the samples, $w^* = \sup_{\mathbf{z}} \pi(\mathbf{z})/q_{\text{def}}(\mathbf{z}; \boldsymbol{\lambda})$, and $C(r, N) > 0$ is a finite constant.

Proof. JSA is described in Algorithm 5. At each iteration, it performs N MCMC transitions and uses the N intermediate states to estimate the gradient. That is,

$$\begin{aligned}
\mathbf{z}_t^{(1)} \mid \mathbf{z}_{t-1}^{(N)}, \boldsymbol{\lambda}_{t-1} &\sim K_{\boldsymbol{\lambda}_{t-1}}(\mathbf{z}_{t-1}^{(N)}, \cdot) \\
\mathbf{z}_t^{(2)} \mid \mathbf{z}_t^{(1)}, \boldsymbol{\lambda}_{t-1} &\sim K_{\boldsymbol{\lambda}_{t-1}}(\mathbf{z}_t^{(1)}, \cdot) \\
&\vdots \\
\mathbf{z}_t^{(N)} \mid \mathbf{z}_t^{(N-1)}, \boldsymbol{\lambda}_{t-1} &\sim K_{\boldsymbol{\lambda}_{t-1}}(\mathbf{z}_t^{(N-1)}, \cdot) \\
\mathbf{g}_{t, \text{JSA}} &= -\frac{1}{N} \sum_{n=1}^N \mathbf{s}(\boldsymbol{\lambda}, \mathbf{z}_t^{(n)}),
\end{aligned}$$

where $K_{\lambda_{t-1}}^n$ is an n -transition IMH kernel using $q_{\text{def}}(\cdot; \boldsymbol{\lambda}_{t-1})$. Under [Assumption 2](#), an IMH kernel is uniformly geometrically ergodic ([Mengersen & Tweedie, 1996](#); [Wang, 2022](#)) as

$$d_{\text{TV}}(K_{\lambda}^k(\mathbf{z}, \cdot), \pi) \leq r^k \quad (7)$$

for any $\mathbf{z} \in \mathcal{Z}$.

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 $\boldsymbol{\eta} = \mathbf{z}^{(1:N)}$ is represented as

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

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

$$\begin{aligned} & d_{\text{TV}}(P_{\lambda}^k(\boldsymbol{\eta}, \cdot), \Pi) \\ &= \sup_A |\Pi(A) - P^k(\boldsymbol{\eta}, A)| \\ &\leq \sup_A \left| \int_A \pi(d\mathbf{z}'^{(1)}) \times \cdots \times \pi(d\mathbf{z}'^{(N)}) \right. \\ &\quad \left. - K_{\lambda}^{(k-1)N+1}(\mathbf{z}^{(1)}, d\mathbf{z}'^{(1)}) \times \cdots \times K_{\lambda}^{kN}(\mathbf{z}^{(N)}, d\mathbf{z}'^{(N)}) \right| \quad \text{Definition of } d_{\text{TV}} \\ &\leq \sup_A \sum_{n=1}^N \left| \int_A \pi(d\mathbf{z}^{(n)}) - K_{\lambda}^{(k-1)N+n}(\mathbf{z}^{(n)}, d\mathbf{z}'^{(n)}) \right| \quad \text{Lemma 1} \\ &= \sum_{n=1}^N d_{\text{TV}}(K_{\lambda}^{(k-1)N+n}(\mathbf{z}^{(n)}, \cdot), \pi) \quad \text{Definition of } d_{\text{TV}} \\ &\leq \sum_{n=1}^N r^{(k-1)N+n} \quad \text{Equation (7)} \\ &= r^{kN} r^{-N} \frac{r - r^{N+1}}{1 - r} \\ &= \frac{r(1 - r^N)}{r^N(1 - r)} (r^N)^k. \end{aligned}$$

Although the constant depends on r 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 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 \(2\)](#).

Furthermore, by [Lemma 2](#), the second moment of the gradient is bounded as

$$\mathbb{E} \left[\|\mathbf{g}_{t, \text{JSA}}\|_2^2 \mid \mathcal{F}_{t-1} \right] = \mathbb{V} \left[\mathbf{g}_{t, \text{JSA}} \mid \mathcal{F}_{t-1} \right] + \text{Bias} \left[\mathbf{g}_{t, \text{JSA}} \mid \mathcal{F}_{t-1} \right]^2 + 2 \text{Bias} \left[\mathbf{g}_{t, \text{JSA}} \mid \mathcal{F}_{t-1} \right] \|\boldsymbol{\mu}\|_2 + \|\boldsymbol{\mu}\|_2^2$$

$$\leq \mathbb{V}[\mathbf{g}_{t,\text{JSA}} | \mathcal{F}_{t-1}] + \text{Bias}[\mathbf{g}_{t,\text{JSA}} | \mathcal{F}_{t-1}]^2 + 2L \text{Bias}[\mathbf{g}_{t,\text{JSA}} | \mathcal{F}_{t-1}] + \|\boldsymbol{\mu}\|_2^2,$$

where $\boldsymbol{\mu} = \mathbb{E}_{\pi} \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z})$. As shown in [Lemma 7](#), the bias terms decreases in a rate of r^{tN} . Therefore,

$$\mathbb{E}[\|\mathbf{g}_{t,\text{JSA}}\|_2^2 | \mathcal{F}_{t-1}] \leq \mathbb{V}[\mathbf{g}_{t,\text{JSA}} | \mathcal{F}_{t-1}] + \|\boldsymbol{\mu}\|_2^2 + \mathcal{O}(r^{tN}). \quad \text{Lemma 7}$$

Note that it is possible to obtain a tighter bound on the bias terms such that $\mathcal{O}(r^{tN}/N)$, if we directly use (K, \mathbf{z}) to bound the bias instead of the higher-level $(P, \boldsymbol{\eta})$ abstraction. The extra looseness comes from the use of [Lemma 1](#).

For the variance term, we show that

$$\begin{aligned} & \mathbb{V}[\mathbf{g}_{t,\text{JSA}} | \mathcal{F}_{t-1}] \\ &= \mathbb{V}\left[\frac{1}{N} \sum_{n=1}^N \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(n)}) \middle| \mathcal{F}_{t-1}\right] \\ &= \frac{1}{N^2} \sum_{n=1}^N \mathbb{V}[\mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(n)}) | \mathcal{F}_{t-1}] + \frac{2}{N^2} \sum_{n=2}^N \sum_{m=1}^{n-1} \text{Cov}(\mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(n)}), \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(m)}) | \mathcal{F}_{t-1}) \\ &\leq \frac{1}{N^2} \sum_{n=1}^N \mathbb{E}\left[\|\mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(n)})\|_2^2 \middle| \mathcal{F}_{t-1}\right] + C_{\text{cov}} \\ &= \frac{1}{N^2} \sum_{n=1}^N \mathbb{E}\left[\|\mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(n)})\|_2^2 \middle| \mathbf{z}_{t-1}^{(N)}, \boldsymbol{\lambda}_{t-1}\right] + C_{\text{cov}} \\ &= \frac{1}{N^2} \sum_{n=1}^N \mathbb{E}_{\mathbf{z}_t^{(n)} \sim K_{\boldsymbol{\lambda}_{t-1}}^n(\mathbf{z}_{t-1}, \cdot)}\left[\|\mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(n)})\|_2^2 \middle| \mathbf{z}_{t-1}^{(N)}, \boldsymbol{\lambda}_{t-1}\right] + C_{\text{cov}} \\ &\leq \frac{1}{N^2} \sum_{n=1}^N \left[n r^{n-1} \mathbb{E}_{\mathbf{z} \sim q_{\text{def}}(\cdot; \boldsymbol{\lambda}_{t-1})}[\|\mathbf{s}(\boldsymbol{\lambda}; \mathbf{z})\|_2^2] + r^n \|\mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_{t-1}^{(N)})\|_2^2 \right] + C_{\text{cov}} \quad \text{Lemma 6} \\ &\leq \frac{1}{N^2} \sum_{n=1}^N [n r^{n-1} L^2 + r^n L^2] + C_{\text{cov}} \quad \text{Assumption 3} \\ &= \frac{L^2}{N^2} \sum_{n=1}^N \left[n \left(1 - \frac{1}{w^*}\right)^{n-1} + \left(1 - \frac{1}{w^*}\right)^n \right] + C_{\text{cov}} \\ &= \frac{L^2}{N^2} \left[(w^*)^2 + w^* - \left(1 - \frac{1}{w^*}\right)^N ((w^*)^2 + w^* + N w^*) \right] + C_{\text{cov}} \\ &= \frac{L^2}{N^2} \left[\frac{1}{2} N^2 + \frac{3}{2} N + \mathcal{O}(1/w^*) \right] + C_{\text{cov}} \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] + C_{\text{cov}}, \end{aligned}$$

where

$$C_{\text{cov}} = \frac{2}{N^2} \sum_{n=2}^N \sum_{m=1}^{n-1} \text{Cov}(\mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(n)}), \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(m)}) | \mathbf{z}_{t-1}^{(N)}, \boldsymbol{\lambda}_{t-1}).$$

The Laurent approximation becomes exact as $w^* \rightarrow \infty$, which is useful considering [Proposition 2](#). \square

Theorem 4. pMCSA, our proposed scheme, is obtained by setting

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

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

$$d_{\text{TV}}(P_{\lambda}^k(\boldsymbol{\eta}, \cdot), \Pi) \leq C(N) r^k \quad \text{and} \quad \mathbb{E}[\|\mathbf{g}_{t,\text{pMCSA}}\|_2^2 | \mathcal{F}_{t-1}] \leq L^2 \left[\frac{1}{N} + \frac{1}{N} \left(1 - \frac{1}{w^*}\right) \right] + \mathcal{O}(r^t) + \|\boldsymbol{\mu}\|_2^2,$$

where $\boldsymbol{\mu} = \mathbb{E}_{\pi} \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z})$, $w^* = \sup_{\mathbf{z}} \pi(\mathbf{z})/q_{\text{def}}(\mathbf{z}; \boldsymbol{\lambda})$ and $C(N) > 0$ is a finite constant.

Proof. Our proposed scheme, pMCSA, is described in [Algorithm 5](#). At each iteration, our scheme performs a single MCMC transition for each of the N samples, or chains, to estimate the gradient. That is,

$$\begin{aligned} \mathbf{z}_t^{(1)} \mid \mathbf{z}_{t-1}^{(1)}, \boldsymbol{\lambda}_{t-1} &\sim K_{\boldsymbol{\lambda}_{t-1}}(\mathbf{z}_{t-1}^{(1)}, \cdot) \\ \mathbf{z}_t^{(2)} \mid \mathbf{z}_{t-1}^{(2)}, \boldsymbol{\lambda}_{t-1} &\sim K_{\boldsymbol{\lambda}_{t-1}}(\mathbf{z}_{t-1}^{(2)}, \cdot) \\ &\vdots \\ \mathbf{z}_t^{(N)} \mid \mathbf{z}_{t-1}^{(N)}, \boldsymbol{\lambda}_{t-1} &\sim K_{\boldsymbol{\lambda}_{t-1}}(\mathbf{z}_{t-1}^{(N)}, \cdot) \\ \mathbf{g}_{t,\text{pMCSA}} &= -\frac{1}{N} \sum_{n=1}^N \mathbf{s}(\boldsymbol{\lambda}, \mathbf{z}_t^{(n)}), \end{aligned}$$

where $K_{\boldsymbol{\lambda}_{t-1}}^n$ is an n -transition IMH kernel using $q_{\text{def}}(\cdot; \boldsymbol{\lambda}_{t-1})$.

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

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

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

$$\begin{aligned} &d_{\text{TV}}(K_\lambda^k(\boldsymbol{\eta}, \cdot), \Pi) \\ &= \sup_A |\Pi(A) - P_\lambda^k(\boldsymbol{\eta}, A)| && \text{Definition of } d_{\text{TV}} \\ &\leq \sup_A \left| \int_A \pi(d\mathbf{z}'_1) \cdot \dots \cdot \pi(d\mathbf{z}'_N) \right. \\ &\quad \left. - K_\lambda^k(\mathbf{z}_1, d\mathbf{z}'_1) \cdot \dots \cdot K_\lambda^k(\mathbf{z}_N, d\mathbf{z}'_N) \right| \\ &\leq \sup_A \sum_{n=1}^N \left| \int_A \pi(d\mathbf{z}'_k) - K_\lambda^k(\mathbf{z}_n, d\mathbf{z}'_n) \right| && \text{Lemma 1} \\ &= \sum_{n=1}^N d_{\text{TV}}(K_\lambda^k(\mathbf{z}_n, \cdot), \pi) && \text{Equation (7)} \\ &\leq \sum_{n=1}^N r^k && \text{Geometric ergodicity} \\ &= N r^k. \end{aligned}$$

Bound on the Gradient Variance By [Lemma 2](#), the second moment of the gradient is bounded as

$$\begin{aligned} \mathbb{E} \left[\|\mathbf{g}_{t,\text{pMCSA}}\|_2^2 \mid \mathcal{F}_{t-1} \right] &= \mathbb{V} \left[\mathbf{g}_{t,\text{pMCSA}} \mid \mathcal{F}_{t-1} \right] + \text{Bias} \left[\mathbf{g}_{t,\text{pMCSA}} \mid \mathcal{F}_{t-1} \right]^2 + 2 \text{Bias} \left[\mathbf{g}_{t,\text{pMCSA}} \mid \mathcal{F}_{t-1} \right] \|\boldsymbol{\mu}\|_2 + \|\boldsymbol{\mu}\|_2^2 \\ &\leq \mathbb{V} \left[\mathbf{g}_{t,\text{pMCSA}} \mid \mathcal{F}_{t-1} \right] + \text{Bias} \left[\mathbf{g}_{t,\text{pMCSA}} \mid \mathcal{F}_{t-1} \right]^2 + 2L \text{Bias} \left[\mathbf{g}_{t,\text{pMCSA}} \mid \mathcal{F}_{t-1} \right] + \|\boldsymbol{\mu}\|_2^2, \end{aligned}$$

where $\boldsymbol{\mu} = \mathbb{E}_\pi \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z})$. As shown in [Lemma 7](#), the bias terms decreases in a rate of r^t . Therefore,

$$\mathbb{E} \left[\|\mathbf{g}_{t,\text{pMCSA}}\|_2^2 \mid \mathcal{F}_{t-1} \right] \leq \mathbb{V} \left[\mathbf{g}_{t,\text{pMCSA}} \mid \mathcal{F}_{t-1} \right] + \|\boldsymbol{\mu}\|_2^2 + \mathcal{O}(r^t). \quad \text{Lemma 7}$$

As noted in the proof of [Theorem 3](#), it is possible to obtain a tighter bound on the bias terms such that $\mathcal{O}(r^t/N)$.

The variance term is bounded as

$$\begin{aligned} &\mathbb{V} \left[\mathbf{g}_{t,\text{pMCSA}} \mid \mathcal{F}_{t-1} \right] \\ &= \mathbb{V} \left[\frac{1}{N} \sum_{n=1}^N \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(n)}) \mid \mathbf{z}_{t-1}^{(1:N)}, \boldsymbol{\lambda}_{t-1} \right] \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{N^2} \sum_{n=1}^N \mathbb{V} \left[\mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(n)}) \mid \mathbf{z}_{t-1}^{(1:N)}, \boldsymbol{\lambda}_{t-1} \right] && \mathbf{z}_t^{(i)} \perp \mathbf{z}_t^{(j)} \text{ for } i \neq j \\
&\leq \frac{1}{N^2} \sum_{n=1}^N \mathbb{E} \left[\left\| \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(n)}) \right\|_2^2 \mid \mathbf{z}_{t-1}^{(1:N)}, \boldsymbol{\lambda}_{t-1} \right]_2 \\
&= \frac{1}{N^2} \sum_{n=1}^N \mathbb{E}_{\mathbf{z}_t^{(n)} \sim K_{\boldsymbol{\lambda}_{t-1}}(\mathbf{z}_{t-1}^{(n)}, \cdot)} \left[\left\| \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_t^{(n)}) \right\|_2^2 \mid \mathbf{z}_{t-1}^{(1:N)}, \boldsymbol{\lambda}_{t-1} \right] \\
&\leq \frac{1}{N^2} \sum_{n=1}^N \left[\mathbb{E}_{\mathbf{z} \sim q_{\text{def}}(\cdot; \boldsymbol{\lambda}_{t-1})} \left[\left\| \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}) \right\|_2^2 \right] + r \left\| \mathbf{s}(\boldsymbol{\lambda}; \mathbf{z}_{t-1}^{(n)}) \right\|_2^2 \right] && \text{Lemma 6} \\
&\leq \frac{1}{N^2} \sum_{n=1}^N [L^2 + rL^2] && \text{Assumption 3} \\
&= \frac{L^2}{N^2} \sum_{n=1}^N [1 + r] \\
&= L^2 \left[\frac{1}{N} + \frac{1}{N} \left(1 - \frac{1}{w^*}\right) \right].
\end{aligned}$$

□

E Additional Experimental Results

E.1 Bayesian Neural Network Regression

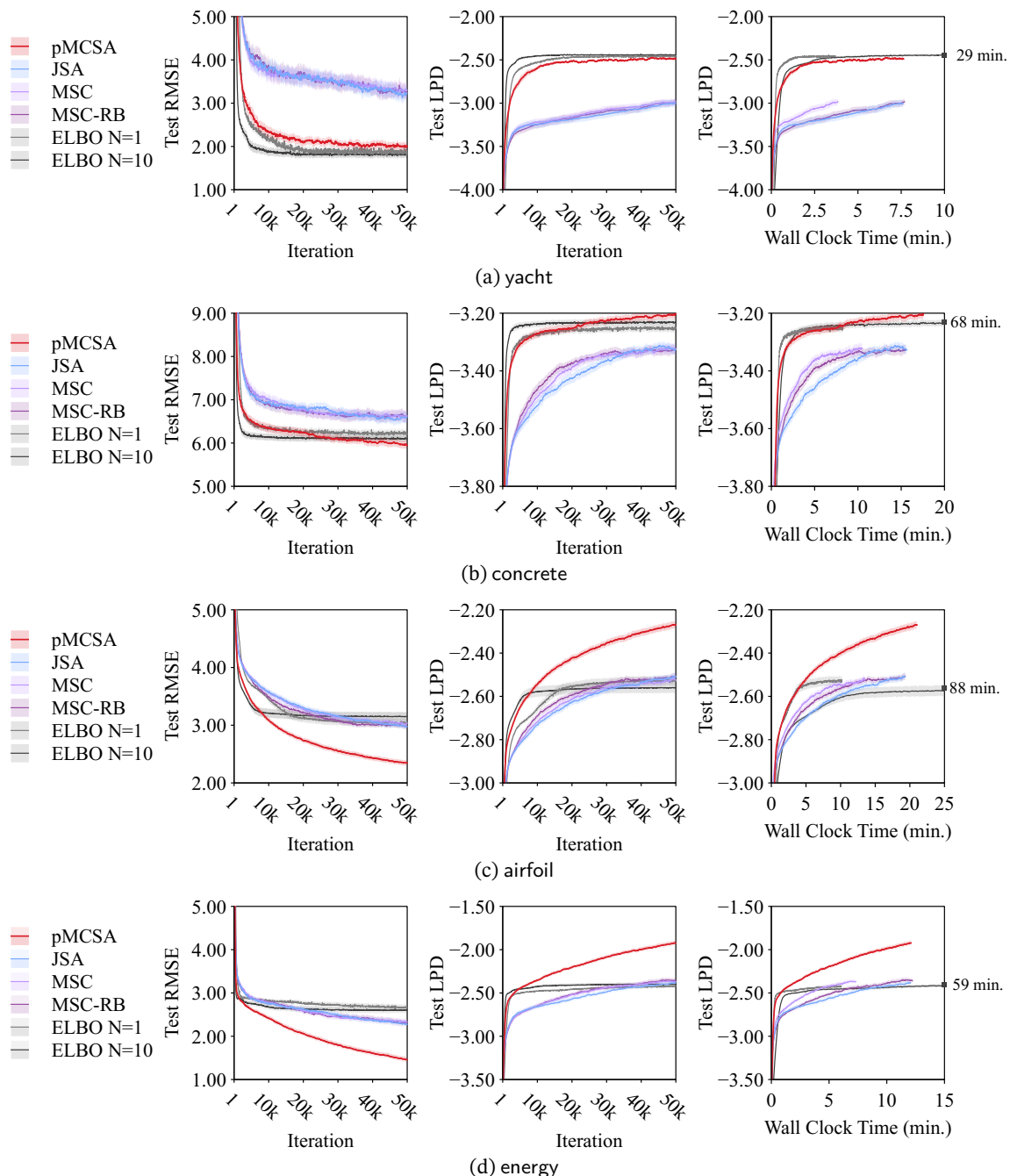


Figure 4: **Test root-mean-square error (RMSE) and test log predictive density (LPD) on Bayesian neural network regression.** The grey squares (■) mark the performance of ELBO $N = 10$ at the wall clock time shown next to it. The error bands show the 95% bootstrap confidence intervals obtained from 20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

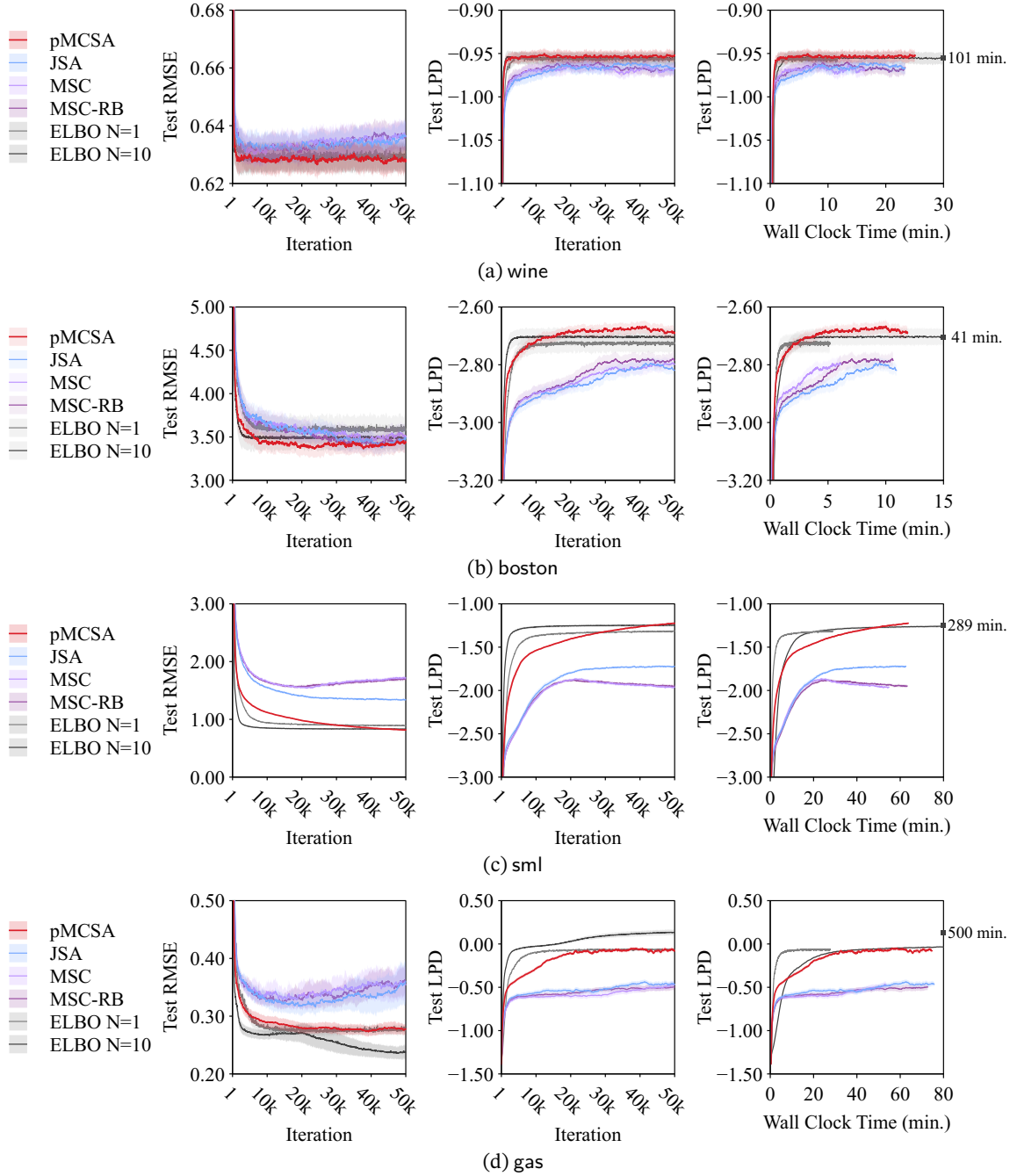


Figure 5: (continued) Test root-mean-square error (RMSE) and test log predictive density (LPD) on Bayesian neural network regression. The grey squares (■) mark the performance of ELBO $N = 10$ at the wall clock time shown next to it. The error bands show the 95% bootstrap confidence intervals obtained from 20 independent 90% train-test splits.

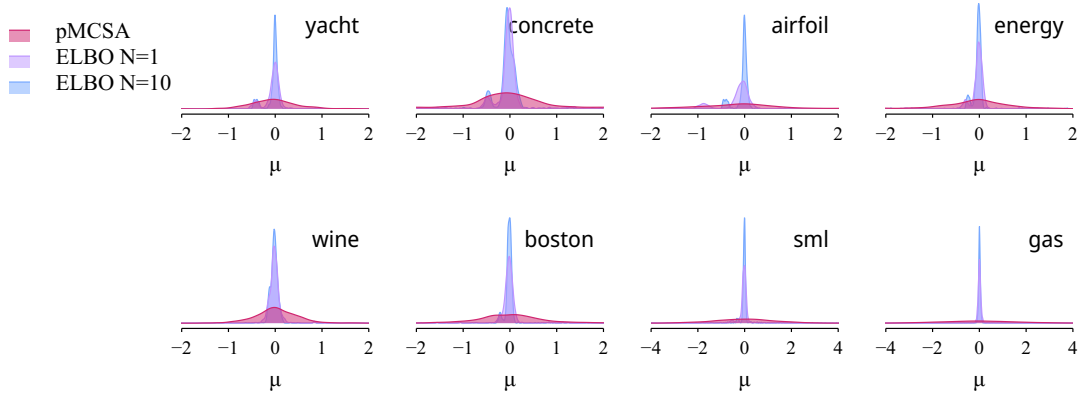


Figure 6: **Distribution of the variational posterior mean of the BNN weights.** The density was estimated with a Gaussian kernel and the bandwidth was selected with Silverman’s rule

E.2 Robust Gaussian Process Regression

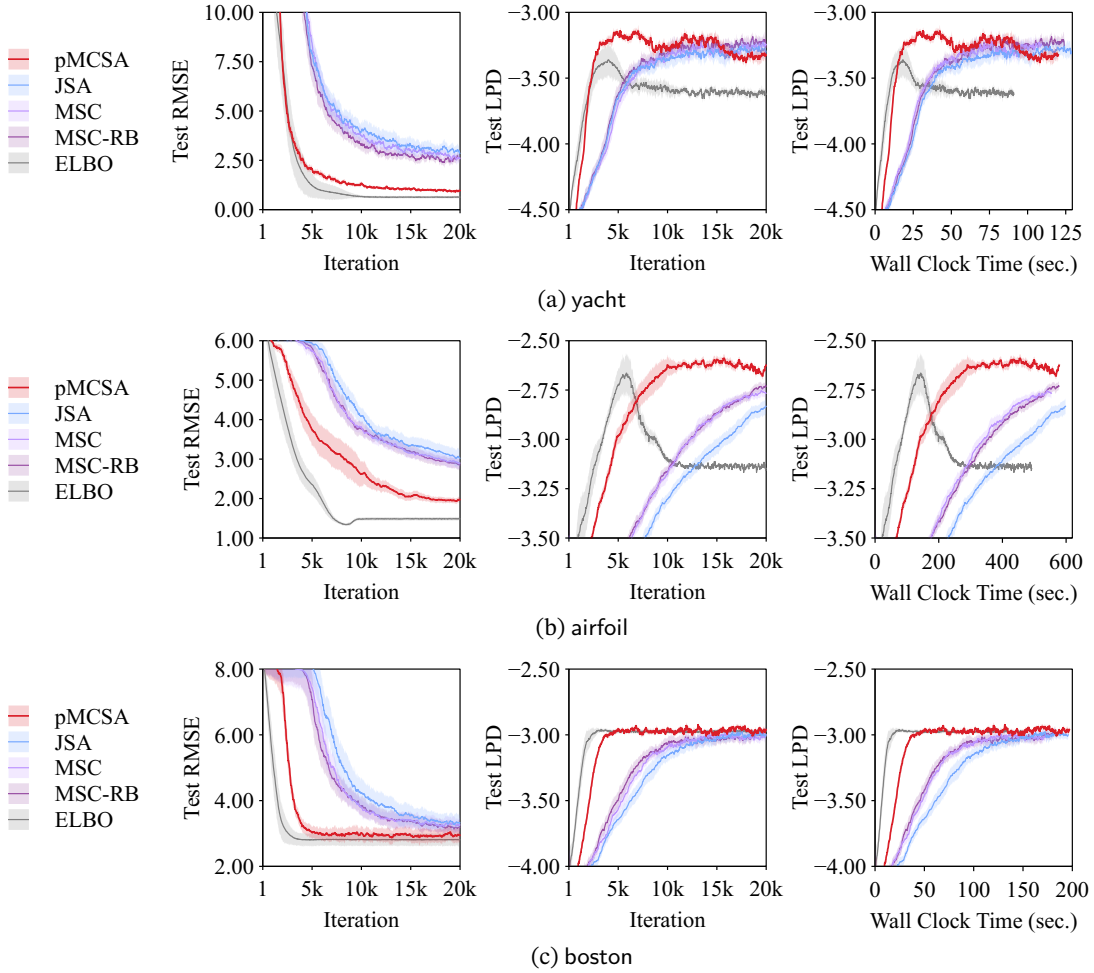


Figure 7: **Test root-mean-square error (RMSE) and test log predictive density (LPD) on robust Gaussian process regression.** The error bands shows the 95% bootstrap confidence interval obtained from 20 repetitions.

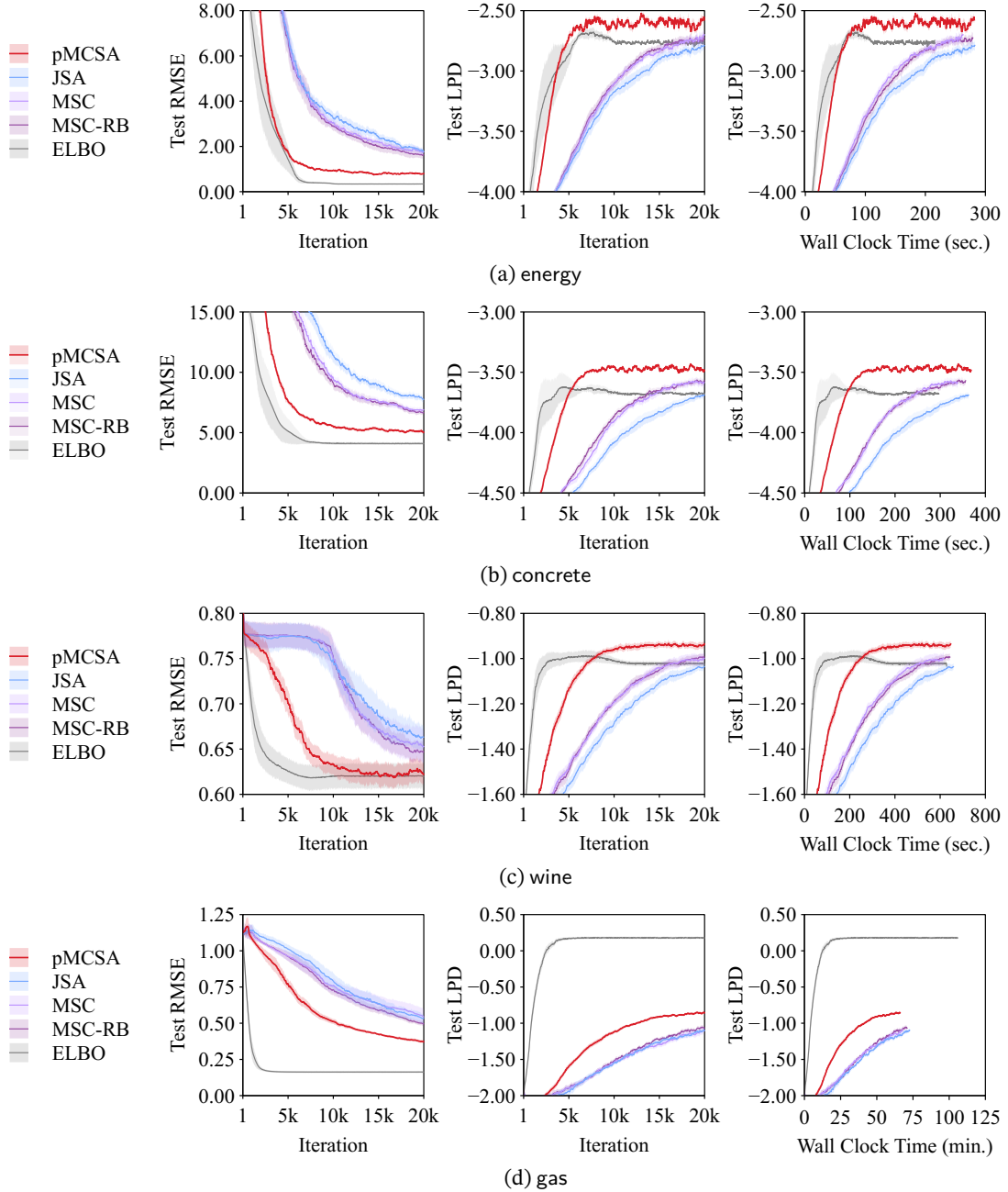


Figure 8: (continued) Test root-mean-square error (RMSE) and test log predictive density (LPD) on robust Gaussian process regression. The error bands show the 95% bootstrap confidence intervals obtained from 20 independent 90% train-test splits.