

Nonlinear MCMC for Bayesian Machine Learning

James Vuckovic
james@jamesvuckovic.com

November 23, 2022

Outline:

1. Introduction
2. Theory
3. Experiments

Motivation — MCMC and Bayesian ML

- ▶ We want to compute the **Bayesian posterior predictive distribution** $P(y|x, \mathcal{D})$.
- ▶ Given a parameteric family of models $P(y, x|\theta)$, we can decompose this problem into the **integral**

$$P(y|x, \mathcal{D}) = \int P(y|x, \theta)P(\theta|\mathcal{D})$$

- ▶ However, this integral is generally **intractable** so we approximate it with samples $\theta^i \sim P(\theta|\mathcal{D})$ i.e.

$$P(y|x, \mathcal{D}) = \int P(y|x, \theta)P(\theta|\mathcal{D}) \approx \frac{1}{N} \sum_{i=1}^N P(y|x, \theta^i)$$

- ▶ However, generating exact samples θ^i is **also intractable** so we can use a **Markov kernel** \mathcal{T} with **invariant measure** $P(\theta|\mathcal{D})$ to generate independent Markov chains $\theta_{n+1}^i \sim \mathcal{T}(\theta_n^i, \bullet)$ and approximate

$$P(y|x, \mathcal{D}) = \int P(y|x, \theta)P(\theta|\mathcal{D}) \approx \frac{1}{N} \sum_{i=1}^N P(y|x, \theta^i) \approx \frac{1}{N} \sum_{i=1}^N P(y|x, \theta_{n_{\text{sim}}}^i)$$

- ▶ This procedure is **Markov Chain Monte Carlo (MCMC)**. The central question of our work: **how do we pick a good transition kernel?**

Nonlinear MCMC — Mean Field System

- ▶ To specify a MCMC algorithm, we need a Markov **transition kernel**
- ▶ Our work directly builds on the paper Andrieu et al., 2011 with some important differences & extensions

On nonlinear Markov chain Monte Carlo

CHRISTOPHE ANDRIEU¹, AJAY JASRA², ARNAUD DOUCET³ and PIERRE DEL MORAL⁴

- ▶ Consider a family of kernels indexed by probability measures $\eta \in \mathcal{P}(\mathbb{R}^d)$

$$K_\eta(x, dy) := (1 - \varepsilon)K(x, dy) + \varepsilon J_\eta(x, dy)$$

- ▶ K is a *linear* Markov kernel, called the **primary** kernel, and
 - ▶ J_η is a family of *nonlinear* “jump **interaction**” Markov kernels
 - ▶ $\varepsilon \in]0, 1[$ is a mixture hyperparameter
- ▶ Let Q be a *linear* Markov kernel called the **auxiliary** kernel
- ▶ We can construct a nonlinear Markov chain $\{(X_n, Y_n)\}_{n=0}^\infty$ from K_η as follows

$$\begin{cases} Y_{n+1} \sim Q(Y_n, \bullet) \\ \eta_{n+1} := \text{Distribution}(Y_{n+1}) & Y_0 \sim \eta_0, X_0 \sim \mu_0 \\ X_{n+1} \sim K_{\eta_{n+1}}(X_n, \bullet) \end{cases} \quad (1)$$

- ▶ We pick Q to be η^\star -invariant, K, J_{η^\star} to be π -invariant

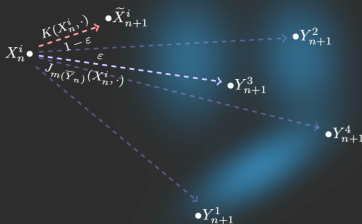
Nonlinear MCMC — Interacting Particle System

- ▶ However equation (1) can't be directly simulated due to $\text{Distribution}(Y_n)$
- ▶ We can *approximate* $\text{Distribution}(Y_n)$ in the mean field system (1) using a set of particles $\bar{Y}_n := \{Y_n^1, \dots, Y_n^N\}$ with the empirical measure

$$\text{Distribution}(Y_n) \approx m(\bar{Y}_n) := \frac{1}{N} \sum_{i=1}^N \delta_{Y_n^i}$$

- ▶ Hence we get the **interacting particle system** which we will *simulate* to obtain MCMC estimates

$$\begin{cases} Y_{n+1}^i \sim Q(Y_n^i, \bullet) \\ \eta_{n+1}^N := m(\bar{Y}_{n+1}) \\ X_{n+1}^i \sim K_N(X_n^i, \bullet) \end{cases} \quad Y_0^i \stackrel{iid}{\sim} \eta_0, \quad X_0^i \stackrel{iid}{\sim} \mu_0, \quad i = 1, \dots, N. \quad (2)$$



Specific Nonlinear Interactions

- ▶ We need to make choices for nonlinear “jump” interactions J_η . We use two proposals in Andrieu et al., 2011

- ▶ Define the potential function $G(x) := \frac{\pi(x)}{\eta^\star(x)}$

- ▶ **Boltzmann-Gibbs transformation**

$$J_\eta^{BG}(x, dy) = \Psi_G(\eta)(dy); \quad \Psi_G(\eta)(dy) = \frac{G(y)}{\eta(G)} \eta(dy)$$

- ▶ This uses G to re-weight the distribution η
- ▶ If $\eta = m(\{Y^1, \dots, Y^N\})$, then this amounts to using a *softmax* over the log-potentials of each particle:

$$\Psi_G(m(\bar{Y})) = \sum_{i=1}^N \frac{G(Y^i)}{\sum_{j=1}^N G(Y^j)} \delta_{Y^i}$$

- ▶ **Accept-Reject Interaction**

$$J_\eta^{AR}(x, dy) = \underbrace{\alpha(x, y)\eta(dy)}_{\text{accept}} + \underbrace{\left(1 - \int \alpha(x, y)\eta(dy)\right)}_{\text{reject}} \delta_x(dy); \quad \alpha(x, y) := 1 \wedge \frac{G(y)}{G(x)}$$

- ▶ This is an “adaptive Metropolis-Hastings” where the proposal distribution is η
- ▶ π is J_{η^\star} -invariant for each of these choices

Theorem 1 (Convergence of Nonlinear MCMC)

Under suitable conditions on K_η and Q , there exist fixed constants $C_1, C_2, C_3 > 0$, a function $\mathcal{R} : [0, \infty[\rightarrow [1, \infty[$, and $\rho > 0$ s.t.

$$\|\mu_n^N - \pi\|_{TV} \leq C_1 \frac{1}{N} \mathcal{R}(1/N) + C_2 \rho^n + C_3 n \rho^n.$$

- ▶ Theorem 1 says it's not *sufficient* to only let $n \rightarrow \infty$ to ensure $\mu_n \rightarrow \pi$
- ▶ It's easy to come up with cases where for any $N > 0$, π is not $K_{m(\bar{Y}_n)}$ -invariant for *any* n . Hence having $N \rightarrow \infty$ is also *necessary* (at least in general)

Corollary 2 (Adapted from Sznitman, 1991, Theorem 2.2)

Suppose that Theorem 1 applies to K_η . Let $\bar{X}_n := \{X_n^1, \dots, X_n^N\}$ be the interacting particle system from (2). Then for every $n \in \mathbb{N}$ and $f \in \mathcal{B}_b(\mathbb{R}^d)$ we have

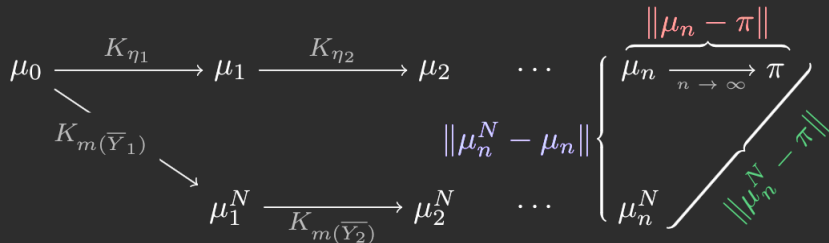
$$\lim_{N \rightarrow \infty} \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N f(X_n^i) - \mu_n(f) \right\| = 0.$$

Proof Sketch — Triangle Inequality

Main ingredients of proof:

1. Ergodicity of K, Q
2. "Regularity" of $\eta \mapsto J_\eta$
3. Use the triangle inequality as follows:

$$\begin{aligned} \|\mu_n^N - \pi\| &\leq \underbrace{\|\mu_n^N - \mu_n\|}_{\text{propagation of chaos}} + \underbrace{\|\mu_n - \pi\|}_{\text{long-time convergence}} \\ &\lesssim \frac{1}{N} \mathcal{R}\left(\frac{1}{N}\right) + \rho^n + n\rho^n \end{aligned}$$

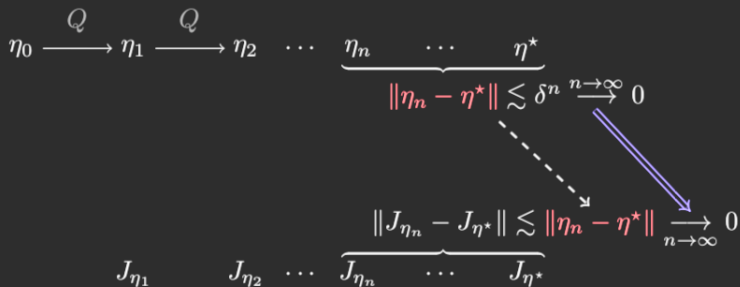


Proof Sketch — Lipschitz Regularity of J

- ▶ We use Lipschitz regularity of J to ensure that as $\eta_n \rightarrow \eta^*$, $J_{\eta_n} \rightarrow J_{\eta^*}$

$$\|J_\eta - J_{\eta'}\| \lesssim \|\eta - \eta'\| \quad \eta, \eta' \in \mathcal{P}(\mathbb{R}^d)$$

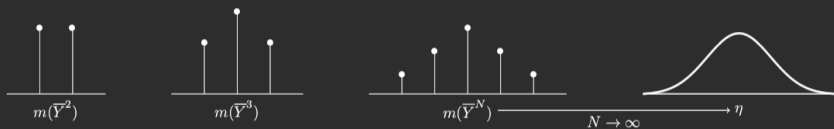
- ▶ This ensures that $K_{\eta_n} \rightarrow K_{\eta^*}$
- ▶ Note: the LHS norm and the RHS norms are *different* – one is over kernels, the other is over probability measures



Proof Sketch — LLN Regularity of J

- ▶ We use particles $\bar{Y} := (Y^1, \dots, Y^N)$ to empirically approximate a measure $\eta \in \mathcal{P}(\mathbb{R}^d)$ by $\eta \approx m(\bar{Y}) = \frac{1}{N} \sum_{i=1}^N \delta_{Y^i}$
- ▶ As $N \rightarrow \infty$, various LLN (or CLT) results tell us $m(\bar{Y}) \rightarrow \eta$
- ▶ We want this same “LLN” regularity of the nonlinear interaction J_η , i.e. as $N \rightarrow \infty$, we want $J_{m(\bar{Y})} \rightarrow J_\eta$
- ▶ Additional considerations:
 - ▶ The convergence should not depend on which $\eta \in \mathcal{P}(\mathbb{R}^d)$ we're approximating
 - ▶ Since \bar{Y} is random $J_{m(\bar{Y})}$ is a *random kernel* which is hard to work with. To simplify, we mean convergence in a suitable *expectation* sense:

$$\mathbb{E}[J_{m(\bar{Y})} f(x)] \rightarrow J_\eta f(x) \quad \forall x \in \mathbb{R}^d, \quad f \text{ in a suitable set of functions}$$

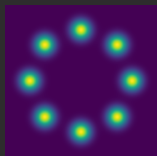
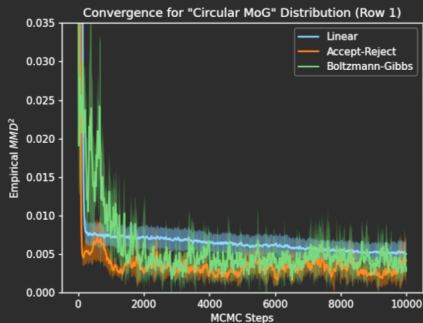


$$\mathbb{E} [J_{m(\bar{Y}^2)}]$$

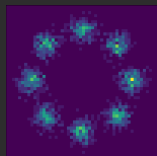
$$\mathbb{E} [J_{m(\bar{Y}^3)}]$$

$$\underbrace{\mathbb{E} [J_{m(\bar{Y}^N)}] - J_\eta}_{\| \mathbb{E} [J_{m(\bar{Y}^N)}] - J_\eta \| \lesssim \frac{1}{N} \mathcal{R}(\frac{1}{N})} \quad J_\eta$$

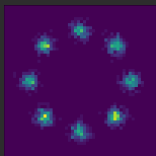
Experiments — 2-dimensional (Circular MoG¹)



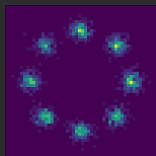
Target
Density



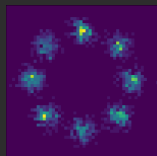
Ground Truth



Linear
(MALA)



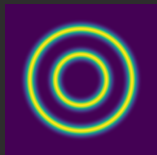
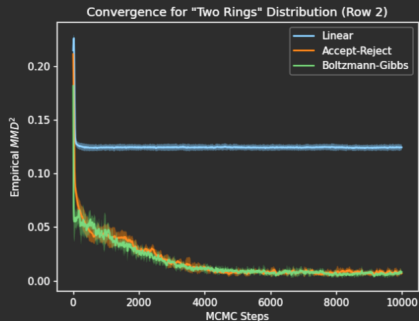
Nonlinear
(AR)



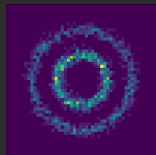
Nonlinear
(BG)

¹Density from Stimper, Schölkopf, and Hernández-Lozano, 2022

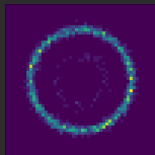
Experiments — 2-dimensional (Two Rings²)



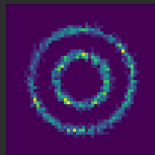
Target Density



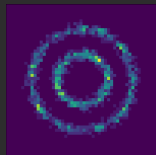
Ground Truth



Linear (MALA)

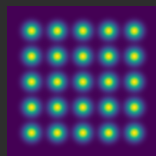
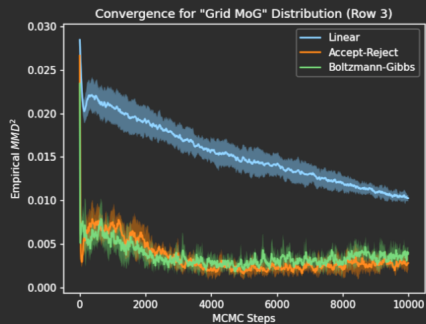


Nonlinear (AR)

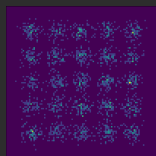


Nonlinear (BG)

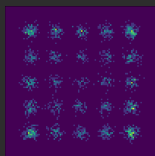
Experiments — 2-dimensional (Grid MoG³)



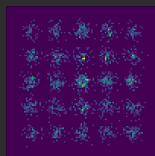
Target
Density



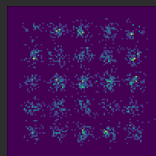
Ground Truth



Linear
(MALA)

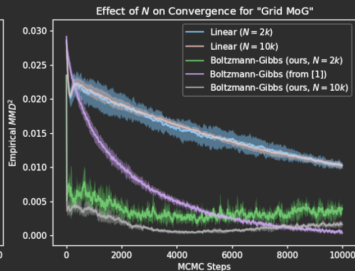
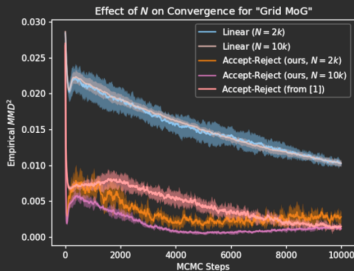
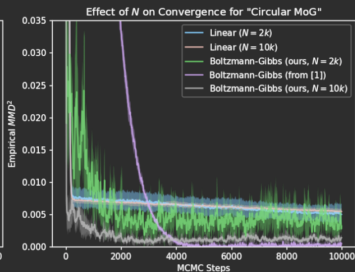
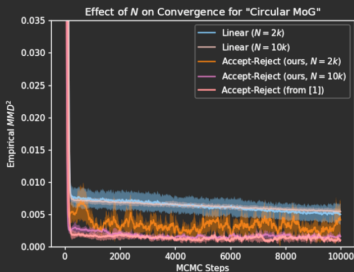


Nonlinear
(AR)



Nonlinear
(BG)

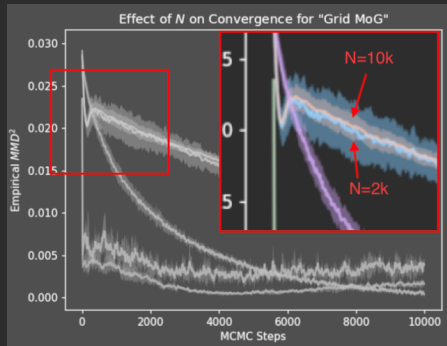
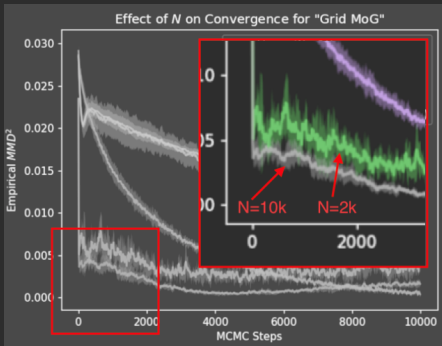
Experiments — 2-dimensional



Accept-Reject

Boltzmann-Gibbs

Experiments — 2-dimensional



- ▶ As predicted by Theorem 1, the nonlinear MCMC convergence rate **depends substantially** on the choice of N (left)
- ▶ This is **not true** for the linear MCMC convergence rate, which only reduces variance by increasing N (right)

Conclusion

► What did we do?






1. We analyzed the convergence of a variation on the family of nonlinear Markov chain Monte Carlo methods proposed in Andrieu et al., 2011
2. Our proof decomposes into two separate results on long-time and large-particle convergence
3. We also applied our theory to two specific choices of samplers also introduced in Andrieu et al., 2011
4. We did some experiments 2-dimensional experiments that demonstrate *superior performance* provided one can choose η^* properly
5. We did some large-scale experiments on CIFAR10 that show our methods are *feasible and comparable but not better* than the linear methods on CIFAR10

► What's next?

1. Investigate how to choose η^* better in high dimension (e.g. for neural networks)
2. Expand to more high-dimensional settings and develop better recipes for MCMC in Bayesian ML

Thank you!

References

-  Andrieu, Christophe et al. (2011). “On nonlinear Markov chain Monte Carlo”. In: *Bernoulli* 17.3, pp. 987–1014. DOI: 10.3150/10-BEJ307. URL: <https://doi.org/10.3150/10-BEJ307>.
-  He, Kaiming et al. (2016). “Deep residual learning for image recognition”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778.
-  Stimper, Vincent, Bernhard Schölkopf, and José Miguel Hernández-Lobato (2022). “Resampling Base Distributions of Normalizing Flows”. In: *International Conference on Artificial Intelligence and Statistics*. PMLR, pp. 4915–4936.
-  Sznitman, Alain-Sol (1991). “Topics in propagation of chaos”. In: *Ecole d’été de probabilités de Saint-Flour XIX—1989*. Springer, pp. 165–251.
-  Zhang, Ruqi et al. (2019). “Cyclical Stochastic Gradient MCMC for Bayesian Deep Learning”. In: *International Conference on Learning Representations*.