

On the Convergence of the Random Walk Metropolis Algorithm

Sam Power

University of Bristol

30 March, 2023

Links & Acknowledgements

- ✂ Main paper today: arXiv 2211.08959;
- ✂ Related: arXiv 2208.05239
- ✂ All joint work with
 - ▶ Christophe Andrieu (Bristol)
 - ▶ Anthony Lee (Bristol)
 - ▶ Andi Q. Wang (Bristol \rightsquigarrow Warwick)
- ✂ Funded by Bayes4Health EPSRC Grant

Talk Goals

- ✿ In recent work, we resolve (to a large extent) the convergence behaviour of the Random Walk Metropolis MCMC algorithm.
- ✿ In this talk, I hope to convey conceptual messages rather than technical details.
- ✿ I want to enable your understanding of how our results can be used to deduce convergence estimates for concrete sampling problems.

Setting: Task

- ✂ Motivating task: making sense of structured probability distributions in high-dimensional spaces
 - ▶ posterior inference in Bayesian statistics
 - ▶ latent variable models, hidden Markov models
 - ▶ generative modeling
 - ▶ non-convex optimisation
 - ▶ ...

Markov Chain Monte Carlo (MCMC)

- ✂ Task: Generate approximate samples from a probability distribution π to which we have *limited access*.
- ✂ MCMC: An iterative approach to this task.
 - ▶ Simulate a time-homogeneous Markov chain $(X_n)_{n \geq 0}$ such that

$$\text{Law}(X_n) \rightarrow \pi \text{ as } n \rightarrow \infty.$$

(and hopefully, quickly)

- ✂ Use samples to ‘understand’ π .

Random Walk Metropolis

✿ Today: Study the *Random Walk Metropolis* (RWM) algorithm

- ▶ Only requires access to density of π , up to a multiplicative constant (typical).
- ▶ Widely-used, simple, 'representative' difficulties

1. At x ,

1.1 Propose $x' \sim \mathcal{N}(x, \sigma^2 \cdot I_d)$.

1.2 Evaluate $r(x, x') = \frac{\pi(x')}{\pi(x)}$.

1.3 With probability $\min\{1, r(x, x')\}$, move to x' ; otherwise, remain at x .

✿ Leaves π invariant; ergodic under mild conditions.

Some Notation

$$\star Q(x, dx') := \mathcal{N}(dx'; x, \sigma^2 \cdot I_d).$$

$$\star \alpha(x, x') := \min \left\{ 1, r(x, x') \right\}.$$

$$\star \alpha(x) := \int Q(x, dx') \alpha(x, x').$$

$$\star \alpha_0 := \inf \left\{ \alpha(x) : x \in \mathbb{R}^d \right\}$$

\star The Random Walk Metropolis kernel P is given by

$$\begin{aligned} P(x, dx') &= Q(x, dx') \cdot \alpha(x, x') \\ &\quad + (1 - \alpha(x)) \cdot \delta(x, dx') \end{aligned}$$

Headlines

✿ Our main results take the (stylised) form

► If $\alpha_0 \gtrsim 1$, then

$$\text{Mixing}(\text{RWM}(\pi, \sigma^2)) \gtrsim \sigma^2 \cdot \text{Mixing}(\text{OLD}(\pi)),$$

where OLD is the continuous-time ‘Overdamped Langevin Diffusion’ process.

✿ Can all be made precise via ‘ L^2 convergence’, ‘mixing times’, etc.

✿ Boils down to

1. Control worst-case acceptance rates (i.e. bound α_0 from below).
2. Understand how ‘nice’ the target π is (i.e. understand mixing of OLD(π)).

The Overdamped Langevin Diffusion

✿ Let $\pi \propto \exp(-U)$ be our target; call U the ‘potential’.

✿ OLD(π) is the SDE given by

$$dX_t = -\nabla U(X_t) dt + \sqrt{2} dW_t.$$

✿ ‘The canonical π -reversible diffusion’.

✿ Fundamental tool for analysing { geometry, concentration, \dots } of π .

✿ Many aspects are well-understood by now.

Crash-Course in Convergence of Langevin Diffusions

1. If U is convex,
 - ▶ then OLD (π) converges at some exponential rate.
2. If U is uniformly quadratically convex,
 - ▶ then OLD (π) initially converges at a *doubly-exponential* rate.
 - ▶ (intuition: 'burn-in' phase is very quick)
3. If U has slower-than linear growth in the tails,
 - ▶ then OLD (π) can only converge at slower-than-exponential rates.
 - ▶ (intuition: time-consuming to get in and out of the tails)

Crash-Course in Convergence of Langevin Diffusions

1. KLS *Conjecture*: if U is convex,
 - ▶ then the exponential rate satisfies

$$\gamma_{\pi} \gtrsim \|\text{Cov}_{\pi}(\text{id})\|_{\text{op}}^{-1}$$

independently of the dimension.

- ▶ (intuition: only bottleneck is the ‘worst’ one-dimensional marginal of π)
2. Various transfer principles (change of measure, transport, . . .).
 - ▶ (intuition: convergence behaviour is robust to various classes of perturbation to π)

Back to the Random Walk Metropolis

- ✂ Rough intuition: when the dimension is high (small moves) and the target is nice (not too rough), the RWM ‘looks like’ OLD.
- ✂ We can show that (again, glossing over details)

$$\text{Mixing}(\text{RWM}(\pi, \sigma^2)) \gtrsim \alpha_0^4 \cdot \sigma^2 \cdot \text{Mixing}(\text{OLD}(\pi)).$$

- ▶ (Actually, this is an instance of a surprisingly general result; can discuss offline.)
- ✂ For RWM, we know that it is not { interesting, relevant } to consider $\alpha_0 \rightarrow 0$.
- ✂ So, how do we control $\alpha_0 \gg 0$?

Regularity Assumptions on Potential

- ✂ Recall that $\alpha(x, x') = \min \left\{ 1, \frac{\pi(x')}{\pi(x)} \right\}$: natural to control $U = -\log \pi$.
- ✂ Smoothness assumption: for some $p \in [1, 2]$, $\psi : \mathbb{R}_+ \rightarrow \mathbb{R}$, it holds that

$$U(x+h) - U(x) - \langle \nabla U(x), h \rangle \leq \psi \left(\|h\|_p \right).$$

- ✂ e.g. ∇U is α -Hölder \rightsquigarrow can take $p = 2$, $\psi(r) \sim r^{1+\alpha}$.
- ✂ $p = 2$ is typically the easiest case to handle.
- ✂ Other p correspond to forms of heterogeneity, roughness.

Acceptance Rate Control for RWM

✿ Lemma: The acceptance rate satisfies

$$\alpha(x) \geq \frac{1}{2} \cdot \exp \left(- \int \mathcal{N}(dz; 0, I_d) \cdot \psi \left(\sigma \cdot \|z\|_p \right) \right),$$

and taking $\sigma = v \cdot d^{-1/p}$ gives that

$$\alpha(x) \geq \frac{1}{2} \cdot \exp \left(-\psi \left(c_p \cdot v \right) + o(1) \right).$$

- ✿ For $p = 2$, consistent with usual optimal scaling results on step-size.
 - n.b. here, we treat *worst-case* rather than *average-case*.
- ✿ For $p < 2$, smaller step-sizes are needed to stabilise α_0 .

Taking Stock

- ✿ So, to ‘know’ the convergence of Random Walk Metropolis on a given target . . . ,
 1. Examine regularity of U , extract a step-size σ which will control α_0 .
 2. Examine global structure of π , characterise the convergence of OLD (π).
 3. ‘Multiply these two things together’.
- ✿ This is the punchline, particularly if you don’t want to get your hands too dirty.

‘Similarity’ of RWM, OLD

- ✿ Morally, the key principle is that $\text{RWM} \approx \text{OLD}$ somehow.
- ✿ In what sense is this true?
- ✿ Not pathwise, nor uniformly (tails).
- ✿ In terms of *exit*, *boundary* behaviour:
 - ▶ If $A \subseteq \mathbb{R}^d$, and $\pi(A)$ is not too tiny, then RWM and OLD both require similar amounts of effort to exit A ; cross the boundary between A , A^c .
 - ▶ Mathematically: isoperimetry, conductance, . . .
- ✿ For L^2 convergence, this is a sufficient ‘resemblance’ between the processes.

Some Examples

✂ Throughout,

- ▶ Assume that initial distance to stationarity is $\exp \Omega(d)$ (typical; achievable).
- ▶ Drop log factors (forgive me).

✂ Model problems:

- ▶ $U(x) = \|x\|_2^2$: take $\sigma \sim d^{-1/2}$, gives $T_{\text{mix}} \lesssim d$
- ▶ $U(x) = \|x\|_2^\alpha$: take $\sigma \sim d^{-1/2}$, gives $T_{\text{mix}} \lesssim d^{2/\alpha}$ for $\alpha \in [1, 2)$.
- ▶ $U(x) = \|x\|_p^\alpha$: take $\sigma \sim d^{-1/p}$, gives $T_{\text{mix}} \lesssim d^{2/p+2/\alpha-1}$ for $\alpha \in [1, 2)$, $p \in [1, 2]$.
 - ▶ One factor for roughness, one factor for tails.
 - ▶ One factor for step-size, one factor for diffusion.

✂ Main example in paper:

- ▶ Suppose that U is convex, with $\text{eigs}(U'') \in [m, L]$ ('well-conditioned').
- ▶ Take $\sigma \sim (L \cdot d)^{-1/2}$.
- ▶ This gives $T_{\text{mix}} \lesssim \kappa \cdot d$, where $\kappa = L/m$.

Recap

- ✂ Random Walk Metropolis for MCMC sampling.
- ✂ Analysis reduces to:
 - ▶ Am I accepting my proposed moves? (roughness)
 - ▶ How would the corresponding Langevin diffusion mix? (concentration)
- ✂ Relatively easy to use, gives sharp results in many cases.