

Optimal scaling for Proximal MALA

Francesca R. Crucinio¹ Alain Durmus² Pablo Jiménez³
and Gareth O. Roberts⁴

¹CREST, ENSAE Paris

²CMAP, École Polytechnique

³Sorbonne Université and Université Paris Cité

⁴Department of Statistics, University of Warwick

31 October 2023

Outline

- 1 Markov chain Monte Carlo
- 2 Proximal MCMC
- 3 Optimal scaling for Proximal MALA

Markov chain Monte Carlo

Aim: sample from a probability distribution π on \mathbb{R}^d and approximate expectations w.r.t. π

$$\int f(x)\pi(x)dx$$

Motivation: compute posterior expectations in Bayesian inference

Metropolis-Hastings

Build a Markov chain $(X_k)_{k \geq 0}$ such that, given the current state X_k

- Sample $Y \sim q(X_k, \cdot)$
- with probability

$$\frac{\pi(Y)q(Y, X_k)}{\pi(X_k)q(X_k, Y)} \wedge 1.$$

set $X_{k+1} = Y$, otherwise, set $X_{k+1} = X_k$.

- ▶ **RWM:** $q(X_k, \cdot) = \mathcal{N}(\cdot; X_k, \sigma^2 I_d)$
- ▶ **MALA:** $q(X_k, \cdot) = \mathcal{N}(\cdot; X_k + \frac{\sigma^2}{2} \nabla \log \pi(X_k), \sigma^2 I_d)$

Outline

- 1 Markov chain Monte Carlo
- 2 Proximal MCMC
- 3 Optimal scaling for Proximal MALA

Proximal MCMC: Idea

- ▶ Build a proposal q using a **proximity map**.
- ▶ Main idea introduced in (Pereyra, 2016) then refined in (Durmus et al., 2018).
- ▶ Can be applied to non-differentiable targets π

Proximity map

Proximity map

For g convex, proper and lower semi-continuous and $\lambda > 0$

$$\text{prox}_g^\lambda(\mathbf{x}) := \arg \min_{\mathbf{u} \in \mathbb{R}^d} \left[g(\mathbf{u}) + \frac{\|\mathbf{u} - \mathbf{x}\|^2}{2\lambda} \right].$$

Moves points in the direction of the minimum of g .

Take $\pi(\mathbf{x}) \propto \exp(-g(\mathbf{x}))$. For each $\lambda > 0$ we can build an approximation π_λ of π .

Moreau-Yoshida envelope

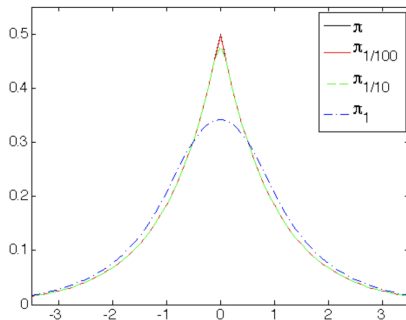


Figure: Moreau-Yoshida envelope for the Laplace distribution $\pi(x) \propto \exp(-|x|)$ (Pereyra, 2016).

Moreau-Yoshida envelope

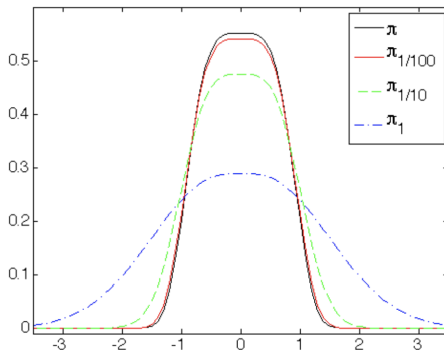


Figure: Moreau-Yoshida envelope for $\pi(x) \propto \exp(-x^4)$ (Pereyra, 2016).

Proximal MCMC

Since π_λ is differentiable for any $\lambda > 0$, we can build a MALA-like algorithm to sample from π_λ :

$$\begin{aligned} q(X_k, Y) &= \mathcal{N} \left(Y; X_k + \frac{\sigma^2}{2} \nabla \log \pi_\lambda(X_k), \sigma^2 I_d \right) \\ &= \mathcal{N} \left(Y; X_k + \frac{\sigma^2}{2\lambda} \left(\text{prox}_g^\lambda(X_k) - X_k \right), \sigma^2 I_d \right), \end{aligned}$$

and then accept/reject using the usual MH step

$$\frac{\pi(Y)q(Y, X_k)}{\pi(X_k)q(X_k, Y)} \wedge 1.$$

Outline

- 1 Markov chain Monte Carlo
- 2 Proximal MCMC
- 3 Optimal scaling for Proximal MALA

Our aim

Investigate the optimal scaling properties of proximal MCMC for targets

$$\pi_d(\mathbf{x}) = \prod_{i=1}^d \pi(x_i) \propto \exp \left(- \sum_{i=1}^d g(x_i) \right)$$

when both $\sigma_d \rightarrow 0$ and $\lambda_d \rightarrow 0$ for

- differentiable g
- the Laplace distribution, $g(x) = |x|$

Notation

We consider

$$\sigma_d^2 = \frac{\ell^2}{d^\alpha}, \quad \lambda_d = \frac{c^2}{2d^\beta}$$

for some $\alpha, \beta > 0$.

Differentiable targets

Given $\pi_d(\mathbf{x}) \propto \exp\left(-\sum_{i=1}^d g(x_i)\right)$, the proposal for proximal MCMC is

$$q(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^d \mathcal{N}\left(y_i; x_i + \frac{\sigma_d^2}{2} g'(\text{prox}_g^{\lambda_d}(x_i)), \sigma_d^2 I_d\right).$$

- $c = 0$ we get $\mathbf{x} + \frac{\sigma_d^2}{2} g'(\text{prox}_g^{\lambda_d}(\mathbf{x})) = \mathbf{x} + \frac{\sigma_d^2}{2} g'(\mathbf{x})$, i.e.
MALA

Differentiable targets – g' Lipschitz

- $\sigma_d^2 = \ell^2 d^{-1/3}$, $\lambda_d = c^2 d^{-\beta}/2$, $\beta > 1/3$
 - ▶ acceptance rate $a(\ell) = 2\Phi\left(-\frac{\ell^3 K_1}{2}\right)$ and $K_1 = K_{MALA}$
 - ▶ speed $h(\ell) = \ell^2 a(\ell)$ maximized at $a(\ell) = 0.574$
- $\sigma_d^2 = \ell^2 d^{-1/3}$, $\lambda_d = c^2 d^{-1/3}/2$
 - ▶ define $r := c^2/\ell^2 > 0$, acceptance rate
 $a(\ell, r) = 2\Phi\left(-\frac{\ell^3 K_2(r)}{2}\right)$ and $K_2^2(r) \geq K_{MALA}^2$ and increasing
 - ▶ speed $h(\ell, r) = \ell^2 a(\ell, r)$ maximized at $a(\ell, r) = 0.574$

Differentiable targets – g' **not** Lipschitz

- $\sigma_d^2 = \ell^2 d^{-1/2}$, $\lambda_d = c^2 d^{-1/4} / 2$
 - ▶ define $r := c^2 / \ell > 0$, acceptance rate $a(\ell, r) = 2\Phi\left(-\frac{\ell^2 K_3(r)}{2}\right)$
and $K_3^2(r) = r^2 \mathbb{E}_X \left[\frac{[g''(X)g'(X)]^2}{4} \right]$
 - ▶ speed $h(\ell, r) = \ell^2 a(\ell, r)$ maximized at $a(\ell) = 0.452$

Differentiable targets – g' Lipschitz

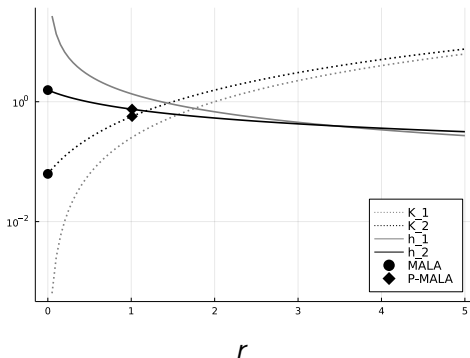


Figure: Speed of Langevin diffusion as a function of $r = c^2 / \ell^2$ for a Gaussian target.

Differentiable targets – g' **not** Lipschitz

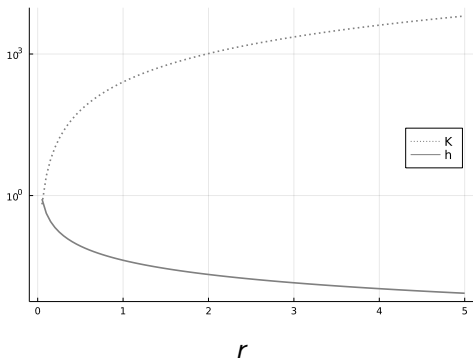


Figure: Speed of Langevin diffusion as a function of $r = c^2/\ell^2$ for a light tail target $\pi(x) \propto \exp(-x^6)$.

Laplace target

Given $\pi_d(\mathbf{x}) \propto \exp\left(-\sum_{i=1}^d |x_i|\right)$, the proposal for proximal MCMC is

$$q(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^d \mathcal{N}(y_i; f(x_i), \sigma_d^2),$$

where

$$f(x_i) = x_i - \frac{\sigma_d^2}{2} \operatorname{sgn}(x_i) \mathbb{1}_{|x_i| \geq \lambda_d}(x_i) - \frac{\sigma_d^2}{2\lambda_d} x_i \mathbb{1}_{|x_i| < \lambda_d}(x_i)$$

Laplace target

- $\sigma_d^2 = \ell^2 d^{-2/3}$, $\lambda_d = c^2 d^{-\beta}/2$, $\beta > 2/3$
 - ▶ acceptance rate $a(\ell) = 2\Phi\left(-\frac{\ell^{3/2}}{(72\pi)^{1/4}}\right)$
 - ▶ speed $h(\ell) = \ell^2 a(\ell)$ maximized at $a(\ell) = 0.360$

The case $c = 0$ corresponds to a subgradient version of MALA.

Take home messages

When implementing proximal MALA we need to take into consideration: efficiency, robustness (i.e. when is the Markov chain geometrically ergodic?), cost of obtaining gradients

We have not explored the robustness of proximal MALA, however,

- if MALA is geometrically ergodic and $\nabla \log \pi(x)$ is cheaper than $\text{prox}_g^\lambda(x) \rightarrow$ use MALA
- in cases in which $\nabla \log \pi(x)$ is more expensive than $\text{prox}_g^\lambda(x) \rightarrow$ use proximal MALA with λ as small as possible
- for light tail distributions use proximal MALA

Thank you!

Bibliography I

Alain Durmus, Eric Moulines, and Marcelo Pereyra. Efficient Bayesian computation by proximal Markov chain Monte Carlo: when Langevin meets Moreau. *SIAM Journal on Imaging Sciences*, 11(1):473–506, 2018.

Marcelo Pereyra. Proximal Markov chain Monte Carlo algorithms. *Statistics and Computing*, 26(4):745–760, 2016.