# Riemann Manifold Langevin and Hamiltonian Monte Carlo Methods

Mark Girolami and Ben Calderhead

Paper reviewed by Hui Li

July 8, 2013

- Standard Metropolis-Hastings Sampling (M-H)
- Metropolis Adjusted Langevin Algorithm (MALA)
- Riemann Manifold Metropolis Adjusted Langevin Algorithm (RMMALA)
- Simplified Manifold Metropolis Adjusted Langevin Algorithm (Simplified mMALA)
- Riemann Manifold Hamiltonian Monte Carlo Method (RMHMC)

Standard Metropolis-Hastings Sampling

- Target distribution  $p(\theta)$
- Proposal distribution  $q(\theta^*|\theta)$
- Draw a new sample sample  $\theta^*$  from proposal distribution  $q(\theta^*|\theta)$
- Accept the new sample with probability  $min\{1, \frac{p(\theta^*)q(\theta|\theta^*)}{p(\theta)q(\theta^*|\theta)}\}$
- Success of MCMC relies upon appropriate proposal design q

- Most-used choice of q(θ\*|θ) includes Normal distribution
   N(θ\*|θ, σ<sup>2</sup>I), a D-dimensional norm distribution with mean θ and covariance matrix σ<sup>2</sup>I
- Random walk
- Need a long time to travel the whole parameter space
- Low acceptance rate
- Poor mixing of the chain and highly-correlated samples
- No information of target distribution needed

Metropolis Adjusted Langevin Algorithm (MALA)

- Stochastic differential equation (SDE) in Langevin diffusion:
   dθ = ∇<sub>θ</sub> L{θ(t)}dt + db(t)
- If *L*{θ(t) is defined as log density of p(θ), then the solution of SDE is the stationary distribution p(θ)
- First order Euler discretization of SDE:  $\theta^* = \theta^n + \varepsilon^2 \nabla_{\theta} \mathcal{L} \{\theta^n\}/2 + \varepsilon \mathbf{z}^n$
- Proposal distribution can be defined as  $\mathcal{N}(\theta^*|\mu(\theta^n, \varepsilon), I)$  with  $\mu(\theta^n, \varepsilon) = \theta^n + \varepsilon^2 \nabla_{\theta} \mathcal{L}\{\theta^n\}/2$

Properties of MALA

- Instead of random walking, MALA tries to sample a new value along the direction which maximizes the log density function
- Isotropic diffusion which forces the choice of step size to accommodate variate with smallest variance
- This can be circumvented by employing a pre-conditioning matrix
   M, N(θ<sup>\*</sup>|μ(θ<sup>n</sup>, ε, M), M<sup>1/2</sup>) with μ(θ<sup>n</sup>, ε, M) = θ<sup>n</sup> + ε<sup>2</sup>M∇<sub>θ</sub>L{θ<sup>n</sup>}/2
- Problem: How to select the pre-conditioning matrix M

A manifold:



- A two-dimensional surface embedded in 3D ambient space;
- Euclidean geometry;
- Non-Euclidean geometry;

How to measure the distance of two points

• Euclidean geometry measures a distance between two points  $\theta = c(t_i)$  and  $\theta + \delta \theta = c(t_i)$  as follows

$$D_E(c(t_i), c(t_j)) = \sqrt{\delta \theta^T \delta \theta}$$
(1)

 Non-Euclidean geometry measures a distance between two points by considering the manifold to which the two points belong

# Riemann manifold

- Tangent space *T<sub>θ</sub>M* is a linear approximation of the manifold and is spanned by the differential operator [*∂*/*∂θ*<sub>1</sub>,..., *∂*/*∂θ<sub>n</sub>*]
- Riemann manifold is a differential manifold in which the tangent space  $T_{\theta}M$  at each point has an inner product defined via a metric tensor  $G_{\theta}$ .
- The metric tensor is a function which takes two vectors  $\theta_1$  and  $\theta_2$  as input and outputs a real-valued scalar  $G_{\theta} : T_{\theta}M \times T_{\theta}M \mapsto \mathcal{R}$
- The distance of two points  $\theta$  and  $\theta + \delta \theta$  in Rieman manifold is defined as  $D_{G_{\theta}}(\theta, \theta + \delta \theta) = \sqrt{\delta \theta^{T} G_{\theta} \delta \theta}$

Properties of Riemannian metric tensor

• Symmetric: 
$$G_{\theta}(t_1, t_2) = G_{\theta}(t_2, t_1);$$

- Bilinear:  $G_{\theta}(t_1 + t_2, t_3) = G_{\theta}(t_1, t_3) + G_{\theta}(t_2, t_3);$
- Positive definite:  $G_{\theta}(t_1, t_1) > 0$

Geometric concept in MCMC

- Denote the expected Fisher information as  $G_{\theta} = cov(\nabla_{\theta}L(\theta))$
- The first order approximation of KL distance between p(y; θ) and p(y; θ + δθ)

$$\begin{aligned} \mathsf{KL}(\theta, \theta + \delta\theta) &= \int \mathsf{p}(\mathsf{y}; \theta + \delta\theta) \log \frac{\mathsf{p}(\mathsf{y}; \theta + \delta\theta)}{\mathsf{p}(\mathsf{y}; \theta)} d\mathsf{y} \\ &\approx \delta\theta^{\mathsf{T}} \mathsf{G}_{\theta} \delta\theta \end{aligned}$$
 (2)

• The expected Fisher information is a Reimannian metric tensor.

Riemannian manifold MALA (RMMALA)

- The SDE defining the Langevin diffusion in Riemann manifold  $d\theta = \tilde{\nabla}_{\theta} \mathcal{L}\{\theta(t)\} dt + d\tilde{\mathbf{b}}(t)$
- The natural gradient  $\tilde{\nabla}_{\theta} \mathcal{L}\{\theta(t)\} = \mathbf{G}^{-1}(\theta(t)) \nabla_{\theta} \mathcal{L}\{\theta(t)\}$
- The Brownian motion in Riemann manifold  $d\tilde{\mathbf{b}}_{i}(t) = |\mathbf{G}^{-1/2}(\theta(t))| + \sum_{j=1}^{D} \frac{\partial}{\partial \theta_{j}} [\mathbf{G}^{-1}(\theta(t))_{ij} \mathbf{G}^{-1/2}(\theta(t))] dt + [\sqrt{\mathbf{G}^{-1}(\theta(t))} d\mathbf{b}(t)]_{i}$
- proposal distribution  $\mathcal{N}(\theta^* | \mu(\theta^n, \varepsilon, \mathbf{G}), \sqrt{\mathbf{G}^{-1}\mathbf{I}})$  with  $\mu(\theta^n, \varepsilon, \mathbf{G}) = \theta^n + \mathbf{G}^{-1}(\theta(t))\nabla_{\theta}\mathcal{L}\{\theta(t)\} + |\mathbf{G}^{-1/2}(\theta(t))| + \sum_{j=1}^{D} \frac{\partial}{\partial \theta_j} [\mathbf{G}^{-1}(\theta(t))_{ij}\mathbf{G}^{-1/2}(\theta(t))] dt$

Illustrative example

- Consider normal density  $p(\mathbf{x}|\mu,\sigma) = N_{\mathbf{x}}(\mu,\sigma)$
- Infer the distribution of parameters  $\mu$  and  $\sigma$  with mMala
- Local inner product on tangent space defined by a metric tensor,
   i.e. δθ<sup>T</sup>G<sub>θ</sub>δθ, where θ = (μ, σ)<sup>T</sup>
- Metric  $G_{\theta}$  is the expected Fisher information matrix

$$G(\mu,\sigma) = \begin{bmatrix} \sigma^{-2} & 0\\ 0 & 2\sigma^{-2} \end{bmatrix}$$
(3)

Metric on tangent space

$$\delta \boldsymbol{\theta}^{\mathsf{T}} \mathbf{G}_{\boldsymbol{\theta}} \delta \boldsymbol{\theta} = \frac{\delta \mu^2 + 2\delta \sigma^2}{\sigma^2} \tag{4}$$

- A sample of size N = 30 was drawn from  $N_x(\mu = 0, \sigma = 10)$
- Starting point is  $\mu_0 = 0, \sigma_0 = 40$



- A sample of size N = 30 was drawn from  $N_x(\mu = 0, \sigma = 10)$
- Starting point is  $\mu_0 = 15, \sigma_0 = 2$



## Conclusions of RMMALA

- Make moves in a Riemann metric rather than according to the standard Euclidean metric
- Utilize the information of the curvature of the manifold

Hamitanian Monte Carlo Method (HMC)

- A joint density  $p(\theta, \mathbf{p})$  is factorized as  $p(\theta, \mathbf{p}) = p(\theta)p(\mathbf{p})$
- $p(\theta)$  is the target density function
- $p(\mathbf{p}) = \mathcal{N}(0, \mathbf{M})$  is an independent auxiliary density function
- The Hamiltanian is defined as the negative of the log joint density  $H(\theta, \mathbf{p}) = -L(\theta) + \frac{1}{2} \log\{(2\pi)^D |\mathbf{M}|\} + \frac{1}{2} \mathbf{p}^T \mathbf{M}^{-1} \mathbf{p}$
- The Hamiltanian equations:

$$\frac{d\theta}{dt} = \frac{\partial H}{\partial \mathbf{p}} = \mathbf{M}^{-1}\mathbf{p}$$
$$\frac{d\mathbf{p}}{dt} = -\frac{\partial H}{\partial \theta} = \nabla_{\theta}L(\theta)$$
(5)

Properties of HMC

- Preserve the total energy:  $H(\theta(t), \mathbf{p}(t)) = H(\theta(0), \mathbf{p}(0))$ , and hence  $p(\theta(t), \mathbf{p}(t)) = p(\theta(0), \mathbf{p}(0))$
- Preserve the volume:  $d\theta(t)d\mathbf{p}(t) = d\theta(0)d\mathbf{p}(0)$
- Time reversal

Euler first-order discretization and Leapfrog method

$$\mathbf{p}(t + \frac{\varepsilon}{2}) = \mathbf{p}(t) + \varepsilon \nabla_{\theta} L(\theta(t))/2$$
  

$$\theta(t + \varepsilon) = \theta(t) + \varepsilon \mathbf{M}^{-1} \mathbf{p}(t + \frac{\varepsilon}{2})$$
  

$$\mathbf{p}(t + \varepsilon) = \mathbf{p}(t + \frac{\varepsilon}{2}) + \varepsilon \nabla_{\theta} L(\theta(t + \varepsilon))/2$$
(6)

 Probability of accepting a new sample (θ\*, p\*) is min{1, exp(-H(θ\*, p\*) + H(θ<sup>n</sup>, p<sup>n+1</sup>))}

### Riemann manifold HMC

- Considering the geometric information, a tensor metric G(θ) defined at a point θ is used instead of M
- Hamiltanian equations become

$$\frac{d\theta}{dt} = \frac{\partial H}{\partial \mathbf{p}} = \mathbf{M}^{-1}\mathbf{p} = \mathbf{G}(\theta)^{-1}\mathbf{p}$$
$$\frac{d\mathbf{p}}{dt} = -\frac{\partial H}{\partial \theta} = \nabla_{\theta}L(\theta) - \frac{1}{2}tr\{\mathbf{G}(\theta)^{-1}\frac{\partial \mathbf{G}(\theta)}{\partial \theta}\}$$
$$+\frac{1}{2}\mathbf{p}^{T}\mathbf{G}(\theta)^{-1}\frac{\partial \mathbf{G}(\theta)}{\partial \theta}\mathbf{G}(\theta)^{-1}\mathbf{p}$$

 Probability of accepting a new sample (θ\*, p\*) is min{1, exp(-H(θ\*, p\*) + H(θ<sup>n</sup>, p<sup>n+1</sup>))} Example – Stochastic volatility model Stochastic volatility model is defined with the latent volatilities taking the form of an AR(1) process such that

$$y_t = \epsilon_t \beta \exp(x_t/2)$$
 (8)

with

$$\mathbf{x}_{t+1} = \phi \mathbf{x}_t + \eta_{t+1} \tag{9}$$

#### where

$$\epsilon_t \sim N(0, 1)$$
  

$$\eta_t \sim N(0, \sigma^2)$$
(10)

and

$$x_1 \sim N(0, \sigma^2/(1-\phi))$$
 (11)

Stochastic volatility model

Joint density

$$p(\mathbf{y}, \mathbf{x}, \beta, \phi, \sigma) = p(\mathbf{x}_1) \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{x}_t, \beta) \prod_{t+2}^T p(\mathbf{x}_t | \mathbf{x}_t - 1, \phi, \sigma)$$
$$\pi(\beta) \pi(\sigma) \pi(\phi)$$
(12)

Split up the sampling procedure in two steps, simulate from

$$egin{aligned} eta, \phi, \sigma | \mathbf{y}, \mathbf{x} &\sim m{p}(eta, \phi, \sigma | \mathbf{y}, \mathbf{x}) \ \mathbf{x} | \mathbf{y}, eta, \phi, \sigma &\sim m{p}(\mathbf{x} | \mathbf{y}, eta, \phi, \sigma) \end{aligned}$$

Metric tensor for parameters

$$G(\beta,\phi,\sigma) = \begin{bmatrix} \frac{2T}{\beta^2} & 0 & 0\\ 0 & 2T & 2\phi\\ 0 & 2\phi & 2\phi^2 + (T-1)(1-\phi^2) \end{bmatrix}$$
(14)

• Metric tensor for latent volatilites  $G(\mathbf{x}) = \frac{1}{2} + \mathbf{C}^{-1}$ 

207

(13)

- The value of  $\beta$  is set to the true value
- The log-joint-probability given different values of σ and φ is shown by the contour plot



# Zoom up the plot



- RMHMC sampling effectively normalizes the gradient in each direction;
- HMC sampling with a unit mass matrix exhibits stronger gradients in horizontal direction than vertical direction and therefore takes much longer to converge to the target density.

Mark Girolami and Ben Calderhead

### Stochastic Volatility Model - Performance

Method	Mean time (s)	$ESS \\ (\beta, \sigma, \phi)$	Standard error ( $\beta, \sigma, \phi$	s/minimum ) ESS	Relative speed
MALA	44.0	(19.1, 11.3, 30.1	) (1.9.0.8.2.1)	3.89	36.7
HMC	424.8	(117, 81, 198)	(9.3, 3.1, 10.3	) 5.24	27.3
MMALA	2455.9	(17.2, 17.4, 44.5	) (2.8, 2.4, 9.2	142.8	1
RMHMC	329.4	(325, 139, 344)	(19.0, 7.3, 25.)	2) 2.37	60.3
Method	Met time	an ESS (1 (s) median,	nînîmum, maxîmum)	s/minimum ESS	Relative speed
					No.
MALA	44	4.0 (9.7, 1)	5.7, 28.4)	4.53	7.5
MALA HMC	44	4.0 (9.7, 1) 4.8 (409, 6	5.7, 28.4) 24, 1239)	4.53	7.5
MALA HMC MMALA	42 424 245	4.0 (9.7, 1) 4.8 (409, 6 5.9 (71.8, 13	5.7, 28.4) 24, 1239) 1.0, 329.8)	4.53 1.04 34.2	7.5 32.9 1

- 2000 simulated observations with  $\beta =$  0.65,  $\sigma =$  0.15 and  $\phi =$  0.98
- 20000 posterior samples averaged over 10 runs.

### Stochastic Volatility Model - Performance



• Posterior marginal density for  $\beta$ ,  $\sigma$  and  $\phi$