

# Riemann manifold Langevin and Hamiltonian Monte Carlo methods

Clément Chadebec

ENS - MVA

January 9, 2020

# Overview

- 1 Rational for new methods
- 2 Hamiltonian Monte Carlo methods
- 3 Parameters Influence
- 4 Riemann Manifold HMC
- 5 Comparison with MCMC Algorithms - Example of Bayesian Logistic Regression
- 6 Conclusion

# Shortcomings of Monte Carlo algorithms

- Not scalable for target densities in high dimension
- Can demonstrate high correlations
- Can demonstrate low acceptance rates  
⇒ Need for new methods

Goal: Simulate a random variable  $\theta \in \mathbb{R}^D \sim \pi$  a target density.

- Introduction of an **independent** auxiliary variable  $\mathbf{Y} \in \mathbb{R}^D \sim \nu = \mathcal{N}(0, M)$  where  $M$  is called the mass matrix
- The negative log-proba of the joint distribution follows:

$$H(\theta, \mathbf{Y}) = - \underbrace{\mathcal{L}(\theta)}_{\text{energy function}} + \frac{1}{2} \log[(2\pi)^D |M|] + \underbrace{\frac{1}{2} \mathbf{Y}^\top M^{-1} \mathbf{Y}}_{\text{kinetic energy}}$$

- The derivatives of  $H$  give

$$\begin{aligned} \frac{d\theta}{d\tau} &= \frac{\partial H}{\partial \mathbf{Y}} = \mathbf{M}^{-1} \mathbf{Y} \\ \frac{d\mathbf{Y}}{d\tau} &= \frac{\partial H}{\partial \theta} = \nabla_{\theta} \mathcal{L}(\mathbf{X}) \end{aligned}$$

- Stormer - Verlet (leapfrog) integrator

$$\mathbf{Y}(\tau + \varepsilon/2) = \mathbf{Y}(\tau) + \varepsilon \nabla_{\theta} \mathcal{L}(\theta)/2$$

$$\theta(\tau + \varepsilon) = \theta(\tau) + \varepsilon M^{-1} \mathbf{Y}(\tau + \varepsilon/2) \quad (\text{Stormer - Verlet})$$

$$\mathbf{Y}(\tau + \varepsilon) = \mathbf{Y}(\tau + \varepsilon/2) + \varepsilon \nabla_{\theta} \mathcal{L}(\theta(\tau + \varepsilon))/2$$

- HMC sampling of  $\pi(\theta)$  as a Gibbs sampler:

$$\mathbf{Y}^{n+1} | \theta^n \sim \mathbf{Y}^{n+1} \sim \mathcal{N}(0, M)$$

$$\theta^{n+1} | \mathbf{Y}^{n+1} \sim \mu(\theta^{n+1} | \mathbf{Y}^{n+1})$$

- $\mu(\theta^{n+1} | \mathbf{Y}^{n+1})$  simulated using Stormer - Verlet scheme and  $(\tilde{\theta}, \tilde{\mathbf{Y}})$  is accepted with probability  $\min\{1, \exp(-H(\tilde{\theta}, \tilde{\mathbf{Y}}) + H(\theta^n, \mathbf{Y}^{n+1}))\}$
- $\implies$  produces an ergodic, time reversible Markov Chain with stationary density  $\pi$
- $\implies$  Difficulty to select  $M$  regardless of the target density

# Leapfrog Impact

Sampling of  $\mathcal{N}(\mathbf{5}, \mathbf{\Sigma})$  where  $\mathbf{\Sigma} = \begin{pmatrix} 0.2 & 0 \\ 0 & 0.8 \end{pmatrix}$

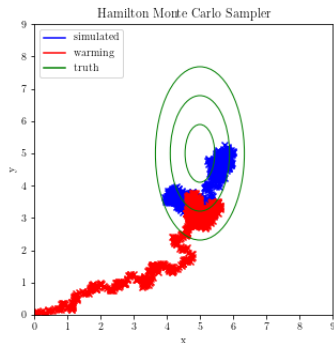


Figure: Leapfrog steps = 2,  $\varepsilon = 0.01$

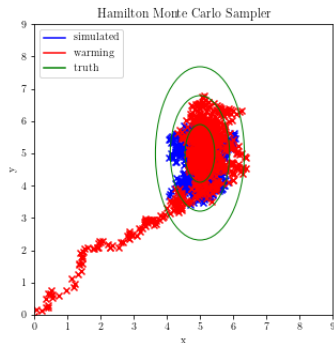


Figure: Leapfrog steps = 5,  $\varepsilon = 0.01$

# Leapfrog Impact

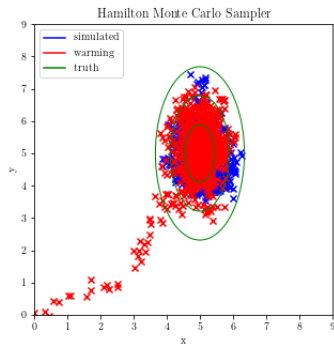


Figure: Leapfrog steps = 10,  $\varepsilon = 0.01$

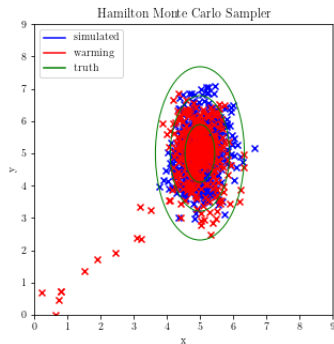


Figure: Leapfrog steps = 20,  $\varepsilon = 0.01$

# Acceptance Ratio

- Dimensions ranging from  $D = 1$  to 50
- Sampling of  $\mathcal{N}(\mathbf{5}, \mathbf{I})$

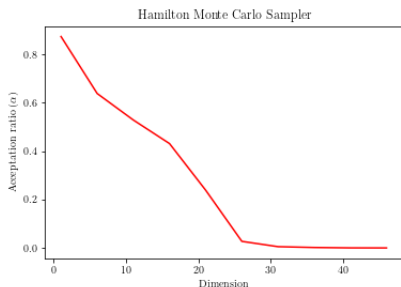


Figure: Leapfrog steps = 20,  $\varepsilon = 0.01$

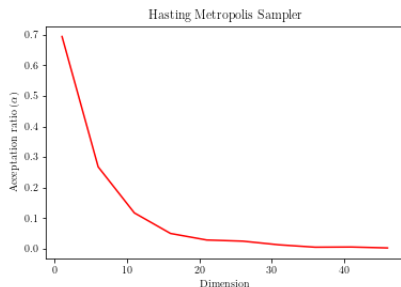


Figure: Sampled from  $\mathcal{N}(\mathbf{0}, \mathbf{I})$



# Auto-correlation

- Sampling of  $\mathcal{N}(\mathbf{5}, \mathbf{I})$
- Leapfrog steps = 20,  $\varepsilon = 0.01$  (HMC)
- Sampled from  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  (HM)

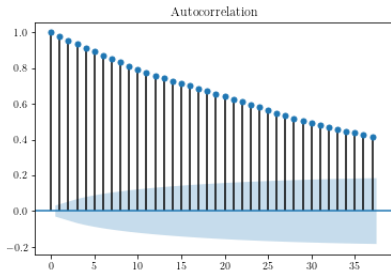


Figure: HMC ( $D = 2$ )

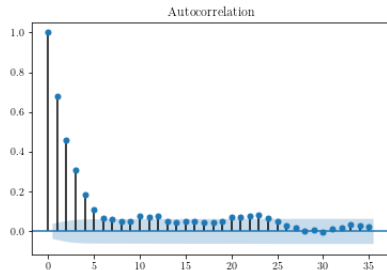


Figure: HM ( $D = 2$ )

# Auto-correlation

- Sampling of  $\mathcal{N}(\mathbf{5}, \mathbf{I})$
- Leapfrog steps = 20,  $\varepsilon = 0.01$  (HMC)
- Sampled from  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  (HM)

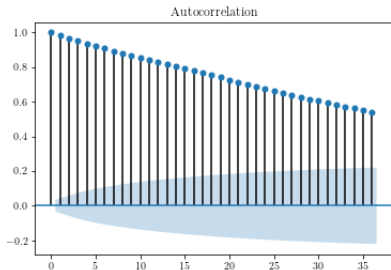


Figure: HMC ( $D = 5$ )

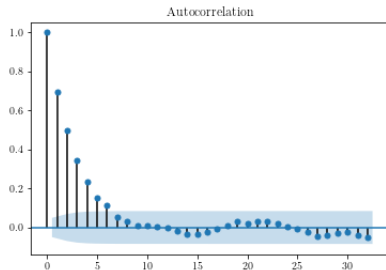


Figure: HM ( $D = 5$ )

# Auto-correlation

- Sampling of  $\mathcal{N}(\mathbf{5}, \mathbf{I})$
- Leapfrog steps = 20,  $\varepsilon = 0.01$  (HMC)
- Sampled from  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  (HM)

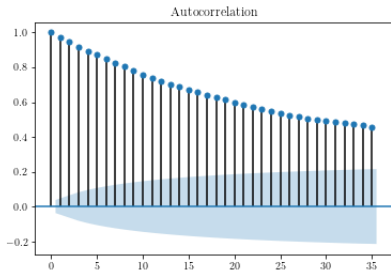


Figure: HMC ( $D = 10$ )

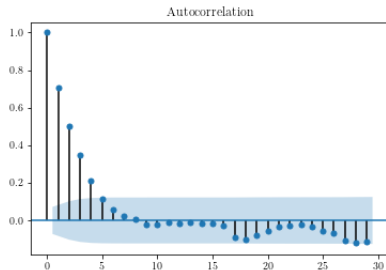


Figure: HM ( $D = 10$ )

# Auto-correlation

- Sampling of  $\mathcal{N}(\mathbf{5}, \mathbf{I})$
- Leapfrog steps = 20,  $\varepsilon = 0.01$  (HMC)
- Sampled from  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  (HM)

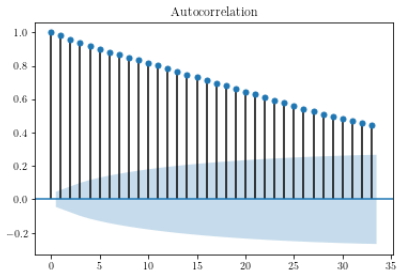


Figure: HMC ( $D = 20$ )

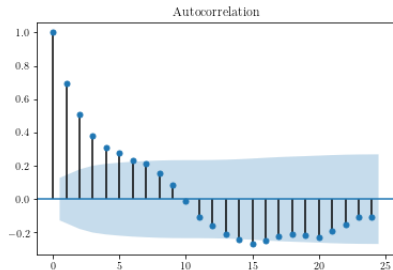


Figure: HM ( $D = 20$ )

# Influence of $M$

- Leapfrog steps = 20,  $\varepsilon = 0.01$
- warm start = 5000

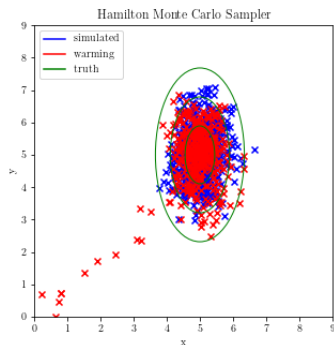


Figure:  $M = 1$

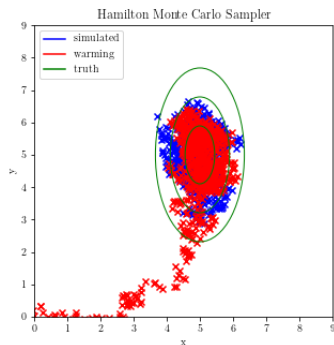


Figure:  $M = 10$

# Influence of $M$

- Leapfrog steps = 20,  $\varepsilon = 0.01$
- warm start = 5000

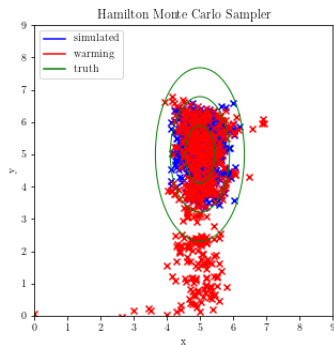


Figure:  $M = \begin{pmatrix} 1 & 0 \\ 0 & 10 \end{pmatrix}$ .

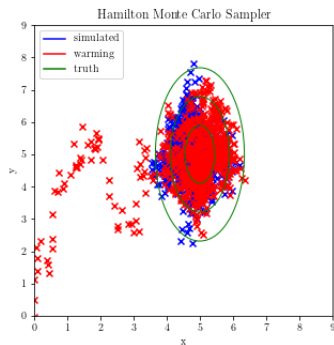


Figure:  $M = \begin{pmatrix} 10 & 0 \\ 0 & 1 \end{pmatrix}$ .

- Distance between two density functions  $\implies p(\mathbf{X}; \theta)$  and  $p(\mathbf{X}; \theta + \delta\theta) = \delta\theta^\top \mathbf{G}(\theta) \delta\theta$  where  $G(\theta)$  is the expected Fisher information matrix

$$\mathbf{G}(\theta) = -\mathbb{E}_{\mathbf{X}|\theta} \left[ \frac{\partial^2}{\partial \theta^2} \log(p(\mathbf{X}|\theta)) \right] = \text{cov} \left[ \frac{\partial}{\partial \theta} \log(p(\mathbf{X}|\theta)) \right]$$

$\implies$  position specific metric on a Riemann manifold

- The Hamiltonian follows

$$H(\theta, \mathbf{Y}) = - \underbrace{\mathcal{L}(\theta)}_{\text{energy function}} + \frac{1}{2} \log[(2\pi)^D |G(\theta)|] + \underbrace{\frac{1}{2} \mathbf{Y}^\top \mathbf{G}(\theta)^{-1} \mathbf{Y}}_{\text{kinetic energy}}$$

- Bayesian approach  $\implies \mathbf{G}(\theta) = -\mathbb{E}_{\mathbf{X}|\theta} \left[ \frac{\partial^2}{\partial \theta^2} \log(p(\mathbf{X}, \theta)) \right]$

- Again the Riemann HMC sampling of  $\pi(\theta)$  can be seen as a Gibbs sampler:

$$\begin{aligned}\mathbf{Y}^{n+1}|\theta^n &\sim \mathcal{N}(0, \mathbf{G}(\theta^n)) \\ \theta^{n+1}|\mathbf{Y}^{n+1} &\sim \mu(\theta^{n+1}|\mathbf{Y}^{n+1})\end{aligned}$$

- $\mu(\theta^{n+1}|\mathbf{Y}^{n+1})$  simulated using generalized Stormer - Verlet scheme and  $(\tilde{\theta}, \tilde{\mathbf{Y}})$  is accepted with probability  $\min\{1, \exp(-H(\tilde{\theta}, \tilde{\mathbf{Y}}) + H(\theta^n, \mathbf{Y}^{n+1}))\}$
- $\implies$  produces an ergodic, time reversible Markov Chain with stationary density  $\pi$
- $M$  mass matrix replaced by position specific metric  $\mathbf{G}(\theta) \longrightarrow$  no need to tune the  $M$  coefficients  
 $\implies$  Need for a new time-reversible numerical integrator for solving the non-separable Hamiltonian
- How to choose the metric  $\mathbf{G}$ ?



# Example of Bayesian Logistic Regression

- The model:
  - Let  $\mathbf{X} \in \mathbb{R}^{N \times D}$  be the design matrix
  - $\beta \in \mathbb{R}^D$  regression parameter with  $\beta \sim \pi = \mathcal{N}(0, \alpha \mathbf{I})$  with  $\alpha$  given
  - We look for  $\beta$  such that  $\mathbf{Y}_n = s(\mathbf{X}_n^\top \beta)$  where  $s$  is the sigmoid function
- The metric tensor follows:

$$\begin{aligned}\mathbf{G}(\beta) &= -\mathbb{E}_{\mathbf{Y}|\beta} \left[ \frac{\partial^2}{\partial \beta^2} \log(p(\mathbf{Y}, \beta)) \right] \\ &= \underbrace{\mathbb{E}_{\mathbf{Y}|\beta} \left[ \frac{\partial^2}{\partial \beta^2} \log(p(\mathbf{Y}|\beta)) \right]}_{\text{Fisher-Rao}} - \underbrace{\frac{\partial^2}{\partial \beta^2} \log(\pi(\beta))}_{\text{Negative Hessian}} \\ &= \mathbf{X}_n^\top \Lambda \mathbf{X}_n + \alpha^{-1} \mathbf{I}\end{aligned}$$

where  $\Lambda$  is diagonal and  $\Lambda_{n,n} = s(\beta^\top \mathbf{X}_n^\top)(1 - s(\beta^\top \mathbf{X}_n^\top))$

# Comparison

- Models considered
  - Component-wise adaptive MH
  - Joint updating Gibbs
  - MALA
  - HMC
  - RHMC - Student
  - Iterated weighted least squares

<i>Name</i>	<i>Covariates (D)</i>	<i>Data points (N)</i>	<i>Dimension of <math>\beta</math> (b)</i>
Pima Indian	7	532	8
Australian credit	14	690	15
German credit	24	1000	25
Heart	13	270	14
Ripley	2	250	7

Figure: Dataset

# Results

- Criteria:  $ESS = N(1 + 2 \sum_k \gamma(k))$  on each covariate.  
 $N$ : the number of posterior samples  
 $\sum_k \gamma(k)$ : sum of the  $K$  monotone sample auto-correlations.

<i>Method</i>	<i>Time</i>	<i>ESS (min, avg, max)</i>	<i>s/min ESS</i>	<i>Relative speed</i>
Metropolis	23.4	(167, 613, 1015)	0.140	13.3
Mala	3.5	(95.5, 316, 667)	0.037	50.3
HMC	117.9	(3182, 3632, 3986)	0.037	50.3
IWLS	7.8	(4.2, 9.9, 69)	1.1862	1
RHMC - S	257.4	(3981, 4934, 5000)	0.065	28.6
<b>RMHMC</b>	<b>246.6</b>	<b>(4757, 5000, 5000)</b>	<b>0.052</b>	<b>35.8</b>

Table: Results ( $D = 24$ ,  $N = 1000$ )

- Strong demonstrated results
- Choice of the metric to be further investigated (Student, ...)
- Choice of the kinetic energy to be further investigated
- What about even bigger dimensions (100, 1000, ...) ?  
⇒ Computation cost scaling

**Proposition:** The transition kernel:

$$P(\theta, A) = \int_{\mathcal{Y}} \mathbf{1}(\tilde{\Phi}^N(\theta, y)) \alpha((\theta, y, \Phi^N(\theta, y))) \nu(y) dy \\ + \mathbf{1}_{\theta}(A) \int_{\mathcal{Y}} (1 - \alpha((\theta, y, \Phi^N(\theta, y))) \nu(y) dy$$

where  $\Phi^N$  is the outcome of  $N$  leapfrog step and  $\tilde{\Phi}(\theta, y) = \theta$

**Proposition:**  $\pi$  is stationary for  $P$

Proof let  $f$  be a Borel function

$$\begin{aligned}
 \mathbb{E}\left[f(\theta^{n+1})|\mathcal{F}_n\right] &= \mathbb{E}_{\mathbf{Y}}\left[\mathbb{E}\left[f(\theta^{n+1})|\mathcal{F}_n, \mathbf{Y}_{n+1}\right]\right] \\
 &= \int_{\mathcal{Y}} \mathbb{E}_{(U, \theta)}\left[f(\tilde{\Phi}^N(\theta, y))\mathbf{I}_{\{U \leq \alpha((\theta, y, \Phi^N(\theta, y)))\}} \right. \\
 &\quad \left. + f(\theta^n)\mathbf{I}_{\{U > \alpha((\theta, y, \Phi^N(\theta, y)))\}}\right] \nu(y) dy \\
 &= \int_{\theta} \int_{\mathcal{Y}} f(\tilde{\Phi}^N(\theta, y)) \alpha((\theta, y, \Phi^N(\theta, y))) \nu(y) dy d\theta \\
 &\quad + \int_{\theta} \delta_{\theta^n}(f) \int_{\mathcal{Y}} (1 - \alpha((\theta, y, \Phi^N(\theta, y)))) \nu(y) dy d\theta
 \end{aligned}$$

Sketch of proof: We use the balanced equation

$$\pi(\theta)P(\theta', \theta) = \pi(\theta')P(\theta, \theta')$$



$$\begin{aligned}\pi(\theta) * \alpha((\theta, y, \theta', y') * \nu(y) &= \pi(\theta) * \min\left(1, \frac{\pi(\theta')\nu(y')}{\pi(\theta)\nu(y)}\right) * \nu(y) \\ &= \pi(\theta') * \alpha(\theta', y', \theta, y) * \nu(y')\end{aligned}$$



$$\begin{aligned}\int_{A \times B} \pi(d\theta)P_2(\theta, d\theta') &= \int_{A \cap B} \pi(d\theta)h(\theta, \cdot) \\ &= \int_{A \cap B} \pi(d\theta')h(\theta, \cdot) \\ &= \int_{A \times B} \pi(d\theta')P_2(\theta', d\theta)\end{aligned}$$