



École des Ponts

ParisTech

Hamiltonian Monte Carlo methods

A Riemannian geometry perspective

Emile Mathieu, Kimia Nadjahi

April 11, 2017

Ecole des Ponts ParisTech

Problem Statement

Problem Statement

Intractable density:

$$p(\beta) = \tilde{p}(\beta) / \int \tilde{p}(\beta) d\beta$$

Metropolis-Hastings:

- Define an ergodic Markov process with stationary distribution $p(\beta)$.
- Transitions $\beta \mapsto \beta^*$ proposed with density $q(\beta^*|\beta)$ accepted with probability

$$\alpha(\beta, \beta^*) = \min \left\{ 1, \frac{\tilde{p}(\beta^*)q(\beta|\beta^*)}{\tilde{p}(\beta)q(\beta^*|\beta)} \right\}$$

What proposal distribution q ?

Typically, random walk: $q(\beta^*|\beta) = \mathcal{N}(\beta^*|\beta, \Lambda)$

Low $\|\Lambda\|$

High acceptance rate

Highly correlated samples

High $\|\Lambda\|$

Low acceptance rate

Not so many correlated samples

Hamiltonian Monte Carlo

Hamiltonian Monte Carlo (I)

Goal:

Make large transitions accepted with high probability.

How to:

- Independent auxiliary variable: $\mathbf{p} \sim \mathcal{N}(\mathbf{p}|\mathbf{0}, \mathbf{M})$
- Joint density: $p(\boldsymbol{\beta}, \mathbf{p}) = p(\boldsymbol{\beta})p(\mathbf{p}) = p(\boldsymbol{\beta})\mathcal{N}(\mathbf{p}|\mathbf{0}, \mathbf{M})$
- The negative joint log-probability is (with $\mathcal{L}(\boldsymbol{\beta}) = \log\{p(\boldsymbol{\beta})\}$)

$$H(\boldsymbol{\beta}, \mathbf{p}) = \underbrace{-\mathcal{L}(\boldsymbol{\beta})}_{\text{potential energy}} + \frac{1}{2} \log\{(2\pi)^D |\mathbf{M}|\} + \underbrace{\frac{1}{2} \mathbf{p}^T \mathbf{M}^{-1} \mathbf{p}}_{\text{kinetic energy}} \quad (1)$$

where $\mathcal{L}(\boldsymbol{\beta}) = \log\{p(\boldsymbol{\beta})\}$

Hamiltonian Monte Carlo (II)

- Hamilton's equations:

$$\frac{d\boldsymbol{\beta}}{d\tau} = \frac{\partial \mathbf{H}}{\partial \mathbf{p}} = \mathbf{M}^{-1}\mathbf{p}, \quad \frac{d\mathbf{p}}{d\tau} = -\frac{\partial \mathbf{H}}{\partial \boldsymbol{\beta}} = \nabla_{\boldsymbol{\beta}}\mathcal{L}(\boldsymbol{\beta}) \quad (2)$$

- Numerical integrator (Stormer-Verlet or leapfrog):

$$\mathbf{p}^{\tau+\frac{\epsilon}{2}} = \mathbf{p}^{\tau} + \frac{\epsilon}{2}\nabla_{\boldsymbol{\beta}}\mathcal{L}\{\boldsymbol{\beta}^{\tau}\} \quad (3)$$

$$\boldsymbol{\beta}^{\tau+\epsilon} = \boldsymbol{\beta}^{\tau} + \epsilon\mathbf{M}^{-1}\mathbf{p}^{\tau+\frac{\epsilon}{2}} \quad (4)$$

$$\mathbf{p}^{\tau+\epsilon} = \mathbf{p}^{\tau+\frac{\epsilon}{2}} + \frac{\epsilon}{2}\nabla_{\boldsymbol{\beta}}\mathcal{L}\{\boldsymbol{\beta}^{\tau+\epsilon}\} \quad (5)$$

- Acceptance probability: $\min(1, \exp\{-H(\boldsymbol{\beta}^*, \mathbf{p}^*) + H(\boldsymbol{\beta}, \mathbf{p})\})$

Hyperparameters:

- Step size ϵ and number of integration steps: via acceptance rate
- Yet, choice of mass matrix \mathbf{M} is critical.

Geometric concepts

Geometric concepts

Goal:

Automatically determine \mathbf{M} at each step.

Fisher-Rao metric:

Distance between parametrized density functions

$$\text{KL}(p(\mathbf{y}; \boldsymbol{\beta}) || p(\mathbf{y}; \boldsymbol{\beta} + \delta\boldsymbol{\beta})) \simeq \delta\boldsymbol{\beta}^T \mathbf{G}(\boldsymbol{\beta}) \delta\boldsymbol{\beta}$$

where $\mathbf{G}(\boldsymbol{\beta}) = -\mathbb{E}_{\mathbf{y}|\boldsymbol{\beta}} \left[\frac{\partial^2}{\partial\boldsymbol{\beta}^2} \log\{p(\mathbf{y}|\boldsymbol{\beta})\} \right]$

Fisher information matrix $\mathbf{G}(\boldsymbol{\beta})$ is p.d. metric defining a Riemann manifold.

General metric tensor

With a Bayesian perspective:

$$\mathbf{G}(\boldsymbol{\beta}) = -\mathbb{E}_{\mathbf{y}|\boldsymbol{\beta}} \left[\frac{\partial^2}{\partial\boldsymbol{\beta}^2} \log\{p(\mathbf{y}, \boldsymbol{\beta})\} \right]$$

Riemann manifold Hamiltonian Monte Carlo

Riemann manifold Hamiltonian Monte Carlo (I)

- The Hamiltonian is

$$\mathbf{H}(\boldsymbol{\beta}, \mathbf{p}) = -\mathcal{L}(\boldsymbol{\beta}) + \frac{1}{2} \log\{(2\pi)^D |\mathbf{G}(\boldsymbol{\beta})|\} + \frac{1}{2} \mathbf{p}^T \mathbf{G}(\boldsymbol{\beta})^{-1} \mathbf{p} \quad (6)$$

- Marginal density:

$$p(\boldsymbol{\beta}) \propto \int \exp\{-\mathbf{H}(\boldsymbol{\beta}, \mathbf{p})\} d\mathbf{p} = \exp\{\mathcal{L}(\boldsymbol{\beta})\}$$

- Joint density: $p(\boldsymbol{\beta}, \mathbf{p}) = p(\boldsymbol{\beta})p(\mathbf{p}|\boldsymbol{\beta}) = p(\boldsymbol{\beta})\mathcal{N}(\mathbf{p}|\mathbf{0}, \mathbf{G}(\boldsymbol{\beta}))$
- Hamilton's equations:

$$\frac{d\beta_i}{d\tau} = \frac{\partial \mathbf{H}}{\partial p_i} = \{\mathbf{G}(\boldsymbol{\beta})^{-1} \mathbf{p}\}_i \quad (7)$$

$$\begin{aligned} \frac{dp_i}{d\tau} = -\frac{\partial \mathbf{H}}{\partial \beta_i} &= \frac{\partial \mathcal{L}(\boldsymbol{\beta})}{\partial \beta_i} - \frac{1}{2} \text{tr} \left\{ \mathbf{G}(\boldsymbol{\beta})^{-1} \frac{\partial \mathbf{G}(\boldsymbol{\beta})}{\partial \beta_i} \right\} \\ &+ \frac{1}{2} \mathbf{p}^T \mathbf{G}(\boldsymbol{\beta})^{-1} \frac{\partial \mathbf{G}(\boldsymbol{\beta})}{\partial \beta_i} \mathbf{G}(\boldsymbol{\beta})^{-1} \mathbf{p} \end{aligned} \quad (8)$$

Riemann manifold Hamiltonian Monte Carlo (II)

Numerical integrator (generalized leapfrog):

$$\mathbf{p}^{\tau+\frac{\epsilon}{2}} = \mathbf{p}^{\tau} - \frac{\epsilon}{2} \nabla_{\beta} \mathbf{H} \left\{ \beta^{\tau}, \mathbf{p}^{\tau+\epsilon/2} \right\} \quad (9)$$

$$\beta^{\tau+\epsilon} = \beta^{\tau} + \frac{\epsilon}{2} \left[\nabla_{\mathbf{p}} \mathbf{H} \left\{ \beta^{\tau}, \mathbf{p}^{\tau+\frac{\epsilon}{2}} \right\} + \nabla_{\mathbf{p}} \mathbf{H} \left\{ \beta^{\tau+\epsilon}, \mathbf{p}^{\tau+\frac{\epsilon}{2}} \right\} \right] \quad (10)$$

$$\mathbf{p}^{\tau+\epsilon} = \mathbf{p}^{\tau+\frac{\epsilon}{2}} - \frac{\epsilon}{2} \nabla_{\beta} \mathbf{H} \left\{ \beta^{\tau+\epsilon}, \mathbf{p}^{\tau+\frac{\epsilon}{2}} \right\} \quad (11)$$

Acceptance probability:

$$\min(1, \exp\{-H(\beta^*, \mathbf{p}^*) + H(\beta, \mathbf{p})\})$$

Application: Bayesian logistic regression

Bayesian logistic regression (BLR)

Probabilistic model

- Likelihood:

$$p(\mathbf{y} = \mathbf{1} | \mathbf{X}) = \eta(\mathbf{X}^T \boldsymbol{\beta})$$

$$\text{with } \eta : x \mapsto (1 + e^{-x})^{-1}$$

- Prior: $\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \alpha \mathbf{I})$
- Metric tensor:

$$\mathbf{G}(\boldsymbol{\beta}) = \mathbf{X}^T \boldsymbol{\Lambda} \mathbf{X} + \alpha^{-1} \mathbf{I}$$

$$\text{with } \Lambda_{i,j} = \eta(\boldsymbol{\beta}^T \mathbf{X}_{n,\cdot}^T) \{1 - \eta(\boldsymbol{\beta}^T \mathbf{X}_{n,\cdot}^T)\}$$

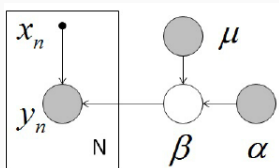


Figure 1: DAG for the Bayesian logistic regression.

Iterative Weighted Least Squares (IWLS)

- MLE of BLR obtained with Newton-Raphson method:

$$\hat{\beta}^{(t)} = ((\alpha I)^{-1} + \mathbf{X}^T \Lambda \mathbf{X})^{-1} (\mathbf{X}^T \Lambda \tilde{\mathbf{y}}(\hat{\beta}^{(t-1)}))$$
$$\hat{\Sigma}_{\beta}^{(t)} = \mathbf{G}(\beta)^{-1} = ((\alpha I)^{-1} + \mathbf{X}^T \Lambda \mathbf{X})^{-1}$$

where $\tilde{\mathbf{y}}(\hat{\beta}^{(t-1)}) = \mathbf{X}\hat{\beta}^{(t-1)} + \Lambda^{-1}(\mathbf{y} - \eta(\beta^T \mathbf{X}))$

- Combination of the MCMC and IWLS iteration schemes:
 1. Initialization: $\beta = \beta^{(0)}$, $t = 1$
 2. Sample β_{new} from $\mathcal{N}(\hat{\beta}^{(t)}, \hat{\Sigma}_{\beta}^{(t)})$
 3. Accept it with probability $\alpha(\beta^{(t-1)}, \beta_{new})$ and set $\beta^{(t)} = \beta_{new}$; otherwise, $\beta^{(t)} = \beta^{(t-1)}$
 4. Do $t := t + 1$ and return to Step 2

Auxiliary variable Gibbs sampler

Representation of BLR with auxiliary variables [Helmes and Hold, 2005]:

$$y_i = \text{sgn}(z_i)$$

$$z_i = \mathbf{X}_i \boldsymbol{\beta} + \epsilon_i$$

$$\epsilon_i \sim \mathcal{N}(0, \lambda_i)$$

$$\lambda_i = (2\psi_i)^2$$

$$\psi_i \sim \text{KS}$$

$$\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \alpha \mathbf{I})$$

$\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\lambda} \sim \mathcal{N}(\mathbf{B}, \mathbf{V})$ with \mathbf{B}, \mathbf{V} as in WLS, $z_i | \boldsymbol{\beta}, \mathbf{X}_i, y_i, \lambda_i \sim$ truncated normal, $\lambda_i | z_i, \boldsymbol{\beta}$ sampled with rejection sampling

Block Gibbs sampler:

1. Update $\{\mathbf{z}, \boldsymbol{\beta}\}$ jointly given $\boldsymbol{\lambda}$
2. Update $\boldsymbol{\lambda} | \mathbf{z}, \boldsymbol{\beta}$

Adaptive Metropolis-Hastings

Component-wise random symmetric walk

At iteration i for component k :

1. Sample $\tilde{\beta}_{i+1}^k \sim \mathcal{N}(\beta_i^k, \sigma_k)$
2. Set $\beta_{i+1}^k = \tilde{\beta}_{i+1}^k$ with probability $\tilde{p}(\tilde{\beta}_{i+1}^k) / \tilde{p}(\beta_i^k)$
3. Otherwise $\beta_{i+1}^k = \beta_i^k$

Component-wise adaptive variance

Every 100 samples after burn-in:

- If $\text{currentAcceptanceRate}^k > \gamma_{max}$
 $\sigma_k = 1.2 * \sigma_k$
- Else if $\text{currentAcceptanceRate}^k < \gamma_{min}$
 $\sigma_k = 0.8 * \sigma_k$

Effective Sample Size (ESS)

$$ESS = \frac{N}{1 + 2 \sum_k \gamma(k)}$$

where N number of samples after burn-in, $\gamma(k)$ autocorrelation of lag k
Ideally, $ESS = N$.

Trade-off time and quality: ratio time / min(ESS)

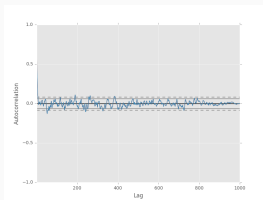
Results

Averaged results: sampling experiments repeated 10 times

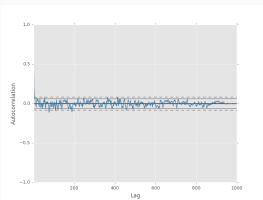
<i>Method</i>	<i>Time</i>	<i>ESS (min, median, max)</i>	<i>time / min(ESS)</i>
Metropolis	10.4	(305.5, 653.6, 801)	0.034
Auxiliary variables	589.4	(710.6, 1199.4, 1715.2)	0.83
HMC	44.9	(3349, 3634, 4141)	0.0134
IWLS	17.9	(21.78, 96.34, 322.4)	0.83
RMHMC	164.1	(4865, 5000, 5000)	0.034

Table 1: Heart data set – comparison of sampling methods

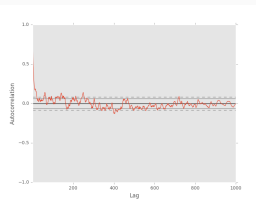
Results



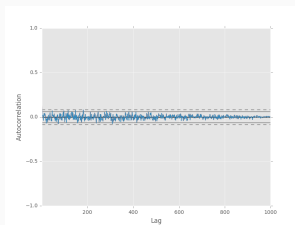
(a) Adaptive MH



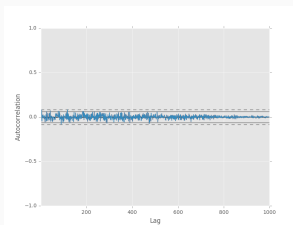
(b) Auxiliary Gibbs



(c) IWLS



(d) HMC



(e) RMHMC

Figure 2: Autocorrelation plots: 1st covariate, 1000 samples of Heart data

Conclusion

Thank you for your attention !

Questions?