

---

# Relativistic Hamiltonian Monte Carlo

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1

## 2 1 Introduction

## 3 2 Relativistic Hamiltonian Monte Carlo

4 Consider a target density  $f(\theta)$  that can be written as  $f(\theta) \propto e^{-U(\theta)}$ . Hamiltonian Monte Carlo  
5 operates by introducing auxiliary variables  $p$  so that  $f(\theta, p) \propto e^{-H(\theta, p)}$ , where

$$H(\theta, p) = U(\theta) + \frac{1}{2m} p^T p \quad (1)$$

6 so that  $p$  is marginally normally distributed. This Hamiltonian lends itself to the interpretation  
7 of a particle with position  $\theta$  and momentum  $p$  moving in a system with potential energy  $U(\theta)$   
8 and according to the classical kinetic energy  $\frac{1}{2m} p^T p$ . We can derive simple update equations for  
9 simulating from these dynamics using Hamilton's equations:

$$\dot{\theta} = \frac{\partial H}{\partial p} \dot{p} = -\frac{\partial H}{\partial \theta} \quad (2)$$

10 giving one possible set of updates (the leapfrog integrator):

$$p_{t+1/2} \leftarrow p_t - \frac{1}{2} \epsilon \nabla U(\theta_t) \quad (3)$$

$$\theta_{t+1} \leftarrow \theta_t + \epsilon \frac{p_{t+1/2}}{m} \quad (4)$$

$$p_{t+1} \leftarrow p_{t+1/2} - \frac{1}{2} \epsilon \nabla U(\theta_{t+1}) \quad (5)$$

11 which is then followed by a Metropolis Hastings accept/reject step. This choice of update is chosen  
12 so that the Hamiltonian  $H$  is left approximately invariant, so that as the acceptance probability  
13 approaches 1. One consequence of these updates is that, when applying HMC to problems where is  
14 very peaked, the momentum  $p$  can become very large, resulting in large updates for  $\theta$ , and thus a very  
15 fine discretization is needed. Consider if, instead of the classical kinetic energy were used for the  
16 Hamiltonian, the relativistic kinetic energy were used instead:

$$K(p) = mc^2 \left( \frac{p^T p}{m^2 c^2} + 1 \right)^{\frac{1}{2}} \quad (6)$$

17 where  $c$  is the "speed of light" which bounds the speed of any particle. This gives the Hamiltonian:

$$H(\theta, p) = U(\theta) + mc^2 \left( \frac{p^T p}{m^2 c^2} + 1 \right)^{\frac{1}{2}} \quad (7)$$

18 The update equations then become

$$p_{t+1/2} \leftarrow p_t - \frac{1}{2}\epsilon \nabla U(\theta_t) \quad (8)$$

$$\theta_{t+1} \leftarrow \theta_t + \epsilon \frac{p_{t+1/2}}{m} \left( \frac{p_{t+1/2}^T p_{t+1/2}}{m^2 c^2} + 1 \right)^{-\frac{1}{2}} \quad (9)$$

$$p_{t+1} \leftarrow p_{t+1/2} - \frac{1}{2}\epsilon \nabla U(\theta_{t+1}) \quad (10)$$

19 Here the momentum is still unbounded and may become very large in the presence of large gradients  
 20 in the potential energy. However, the size of the  $\theta$  update is bounded by  $c$ , thus the behavior of the  
 21 proposed samples can be more easily controlled in the presence of large gradients. The marginal  
 22 distribution of  $p$  is no longer normal, but its density is log-concave and can be sampled using Adaptive  
 23 Rejection Sampling.

### 24 3 Stochastic Gradient Relativistic Hamiltonian Monte Carlo

25 Hamiltonian Monte Carlo algorithms are also of particular interest for “stochastic gradient” style  
 26 algorithms where mini-batches are used to form noisy estimates of the gradients. One motivation for  
 27 this is that the momentum serves as a reservoir of previous gradient information; a large gradient  
 28 will result in a large  $p$ , which may stay large for a while unless met with another large gradient, thus  
 29 retaining the memory of a strong signal on prior batches of data. However, due to potentially large  
 30 variability in the gradient computed in these algorithms, stochastic gradient Hamiltonian algorithms  
 31 may still result in overly large updates, requiring very small values of  $\epsilon$  and thus potentially slow  
 32 convergence. This motivates the use of the Relativistic Hamiltonian in a stochastic gradient sampler;  
 33 the inherent bound in the update size allows the sampler to more easily smooth out the noise in the  
 34 gradient over multiple steps.

35 Ma et al. [2015] gives a framework for taking update equations associated with a particular Hamilto-  
 36 nian and constructing asymptotically consistent stochastic gradient samplers. Specifically, Ma et al.  
 37 [2015] consider a SDE with drift  $f(z)$  and diffusion  $2D(z)$ :

$$dz = f(z)dt + \sqrt{2D(z)}dW_t \quad (11)$$

38 where  $z = (\theta, p)$ ,  $W_t$  is a  $d$ -dimensional Wiener process, and

$$f(z) = -[D(z) + Q(z)] \nabla H(z) + \Gamma(z), \Gamma_i = \sum_{j=1}^d \frac{\partial}{\partial z_j} (D_{ij}(z) + Q_{ij}(z)) \quad (12)$$

39 where  $Q(z)$  is skew-symmetric. Then the update equations

$$z_{t+1} \leftarrow z_t - \epsilon_t \left[ [D(z_t) + Q(z_t)] \nabla \tilde{H}(z_t) + \Gamma(z_t) \right] + \mathcal{N}(0, \epsilon_t(2D(z_t) - \epsilon_t \hat{B}_t)) \quad (13)$$

40 gives an asymptotically consistent chain when the stepsizes  $\epsilon_t$  decrease to zero at the appropriate rate.  
 41 Here  $\tilde{H}(z)$  is the estimate of the Hamiltonian, e.g. using mini-batches, and  $\hat{B}$  is an estimate of the  
 42 variance of the noise of the approximate gradient computation. Note that this estimate need not be  
 43 unbiased for the chain to be consistent – failing better choices we may choose  $\hat{B}_t = 0$ . In practice,  
 44 decreasing the stepsizes  $\epsilon_t$  results in progressively slower mixing, and it is often preferable to fix a  
 45 stepsize and accept that the sampler will incur some asymptotic bias. We can formulate Relativistic  
 46 Hamiltonian Monte Carlo into this framework by taking

$$H(\theta, p) = U(\theta) + mc^2 \left( \frac{p^T p}{m^2 c^2} + 1 \right)^{\frac{1}{2}} \quad (14)$$

$$D(z) = \begin{bmatrix} 0 & 0 \\ 0 & D \end{bmatrix} \quad (15)$$

$$Q(z) = \begin{bmatrix} 0 & -I \\ I & 0 \end{bmatrix} \quad (16)$$

47 which gives  $\Gamma_i(z) = 0$  and

$$f \left( \begin{bmatrix} \theta \\ p \end{bmatrix} \right) = - \begin{bmatrix} 0 & -I \\ I & D \end{bmatrix} \begin{bmatrix} \nabla U(\theta) \\ \frac{p}{m} \left( \frac{p^T p}{m^2 c^2} + 1 \right)^{-\frac{1}{2}} \end{bmatrix} \quad (17)$$

48 which gives the SDE

$$d\theta = \frac{p}{m} \left( \frac{p^T p}{m^2 c^2} + 1 \right)^{-\frac{1}{2}} \quad (18)$$

$$dp = \left( -\nabla U(\theta) - D \frac{p}{m} \left( \frac{p^T p}{m^2 c^2} + 1 \right)^{-\frac{1}{2}} \right) dt + \sqrt{2D} dW_t \quad (19)$$

49 Then (13) gives the updates:

$$p_{t+1} \leftarrow p_t - \epsilon_t \nabla \tilde{U}(\theta_t) - \epsilon_t D \frac{p_t}{m} \left( \frac{p_t^T p_t}{m^2 c^2} + 1 \right)^{-\frac{1}{2}} + \mathcal{N}(0, \epsilon_t (2D - \epsilon_t \hat{B}_t)) \quad (20)$$

$$\theta_{t+1} \leftarrow \theta_t + \epsilon_t \frac{p_{t+1}}{m} \left( \frac{p_{t+1}^T p_{t+1}}{m^2 c^2} + 1 \right)^{-\frac{1}{2}} \quad (21)$$

## 50 4 A Stochastic Gradient Nosé-Hoover Thermostat for Relativistic 51 Hamiltonian Monte Carlo

52 The stochastic gradient version of HMC (SGHMC) introduced by Chen et al. [2014] can be improved  
53 by introducing an additional dynamic variable  $\xi$  to adaptively increase or decrease the momenta  
54 (Ding et al. [2014], Leimkuhler and Shang [2016]). The extended systems has a Hamiltonian of the  
55 form

$$H(\theta, p, \xi) = U(\theta) + \frac{1}{2} p^T p + \frac{d}{2} (\xi - D)^2 \quad (22)$$

56 The dynamics of this approach, known as stochastic gradient Nosé-Hoover thermostat due to its links  
57 to statistical physics, can be expressed as:

$$d\theta = p dt \quad (23)$$

$$dp = -\nabla \tilde{U} dt - \xi p dt + \sqrt{2D} dW_t \quad (24)$$

$$d\xi = \frac{1}{d} (p^T p - 1) dt \quad (25)$$

58 Intuitively this approach works because

$$\mathbb{E} \left[ \frac{d\xi}{dt} \right] = 0, \text{ when sampling from the target joint distribution} \quad (26)$$

59 The system adaptively ‘heats’ or ‘cools’ to push the system closer to obeying (26). Hence the  
60 additional dynamics will move the distribution closer to the equilibrium. In particular this helps to  
61 reduce the bias of SGHMC. A natural question is whether these methods can be extended to relativistic  
62 HMC. Leimkuhler and Shang [2016] show that for a general kinetic energy  $K(p)$ , provided that  $\xi$  is  
63 normally distributed in equilibrium (i.e. using the Hamiltonian in (22)) the  $\xi$  dynamics become

$$d\xi = \frac{1}{d} (\|\nabla K(p)\|^2 - \nabla^2 K(p)) dt \quad (27)$$

64 Note that these general dynamics can still be interpreted as maintaining an equation like (26) since

$$\mathbb{E} \left[ \frac{\partial^2 K}{\partial p_i^2} \right] = \int \frac{\partial^2 K}{\partial p_i^2} e^{-K(p)} dp \quad (28)$$

$$= \underbrace{\int \left[ \frac{\partial K}{\partial p_i} e^{-K(p)} \right]_{p_i=-\infty}^{p_i=\infty} dp}_{=0} - \int \frac{\partial K}{\partial p_i} \left( -\frac{\partial K}{\partial p_i} e^{-K(p)} \right) dp = \mathbb{E} \left[ \left( \frac{\partial K}{\partial p_i} \right)^2 \right] \quad (29)$$

65 and hence  $\mathbb{E} \left[ \frac{d\xi}{dt} \right] = 0$ . We can fit these ideas into the framework of Ma et al. [2015] by defining:

$$H(\theta, p, \xi) = U(\theta) + K(p) + \frac{d}{2} (\xi - D)^2 \quad (30)$$

$$D(\theta, p, \xi) = \begin{pmatrix} 0 & 0 & 0 \\ 0 & D \cdot I & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad (31)$$

$$Q(\theta, p, \xi) = \begin{pmatrix} 0 & -I & 0 \\ I & 0 & \nabla K(p)/d \\ 0 & -\nabla K(p)^T/d & 0 \end{pmatrix} \quad (32)$$

66 This gives

$$\Gamma = \begin{pmatrix} 0 \\ 0 \\ -\nabla^2 K(p)/d \end{pmatrix} \quad (33)$$

67 and the dynamics become

$$d\theta = \nabla K(p)dt \quad (34)$$

$$dp = -\nabla \tilde{U}dt - \xi \nabla K(p)dt + \sqrt{2D}dW_t \quad (35)$$

$$d\xi = \frac{1}{d} (\|\nabla K(p)\|^2 - \nabla^2 K(p)) dt \quad (36)$$

68 This gives a general recipe for a stochastic gradient Nosé-Hoover thermostat with a general kinetic  
69 energy  $K(p)$ . For the relativistic kinetic energy we find

70 Add relativistic dynamics and updates from Xiaoyu's notes

## 71 5 Experiments

## 72 6 Conclusion

## 73 References

- 74 Tianqi Chen, Emily Fox, and Carlos Guestrin. Stochastic Gradient Hamiltonian Monte Carlo. In  
75 *Proceedings of The 31st International Conference on Machine Learning*, pages 1683–1691, 2014.  
76 URL <http://jmlr.org/proceedings/papers/v32/chen14>.
- 77 Nan Ding, Youhan Fang, Ryan Babbush, Changyou Chen, Robert D. Skeel, and Hartmut Neven.  
78 Bayesian Sampling Using Stochastic Gradient Thermostats. In *Advances in Neural Informa-*  
79 *tion Processing Systems*, pages 3203–3211, 2014. URL [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/5592-bayesian-sampling-using-stochastic-gradient-thermostats)  
80 [5592-bayesian-sampling-using-stochastic-gradient-thermostats](http://papers.nips.cc/paper/5592-bayesian-sampling-using-stochastic-gradient-thermostats).
- 81 Benedict Leimkuhler and Xiaocheng Shang. Adaptive Thermostats for Noisy Gradient Systems.  
82 *SIAM Journal on Scientific Computing*, 38(2):A712–A736, mar 2016. ISSN 1064-8275. doi:  
83 [10.1137/15M102318X](https://doi.org/10.1137/15M102318X). URL <http://epubs.siam.org/doi/10.1137/15M102318X>.
- 84 Yi-An Ma, Tianqi Chen, and Emily B. Fox. A Complete Recipe for Stochastic Gradient MCMC. jun  
85 2015. URL <http://arxiv.org/abs/1506.04696>.