

663 Final Project Progress Report

Jake Coleman and Sayan Patra

1 Abstract

For our project we will use Neal’s 2011 paper, “MCMC using Hamiltonian dynamics.” The paper discusses how to use Hamiltonian dynamics as a sampling scheme to explore target spaces better than traditional Metropolis-Hastings algorithms. The Hamiltonian is the sum of potential energy (based on position) and kinetic energy (based on momentum) - Hamilton’s equations relate the two partial derivatives of the Hamiltonian to each other, and define a mapping from the state at time t to the state at time $t + s$. In Hamiltonian Monte Carlo (HMC), we draw auxiliary momentum variables from a Gaussian distribution, and use Hamiltonian dynamics simulations to update the position variable (which follows the distribution of interest). At the end of a user-defined number of steps of simulation, the new variables are accepted or rejected in a Metropolis-Hastings step.

In this report we will explore basic HMC with the “Leapfrog” discretization method, and follow some examples (such as highly-correlated multivariate Gaussian distributions) comparing HMC to random-walk Metropolis Hastings that show improvement for HMC. We will establish the superiority of the HMC method over regular random walk sampling schemes in the case of higher dimensions. We will also implement an extension of HMC proposed by Neal (1994) that uses “windows” of states to allow for a high probability of acceptance for all trajectories. Finally, we plan on converting the code to Cython or JIT to speed up implementation, and compare to existing HMC packages.

2 Introduction

Markov Chain Monte Carlo (MCMC) has been used to simulate distributions since the landmark paper by Metropolis et al. (1953). In this report we follow Neal (2011) in his exploration of Hamiltonian Monte Carlo (HMC), a variant of MCMC which utilizes properties physical properties of molecules known as “*Hamiltonian dynamics*.” Hamiltonian dynamics rely on deterministic motion of molecules described by Newton’s laws of motion, and were only applied as part of an MCMC algorithm in 1987 by Duane, Kennedy, Pendleton, and Roweth, with statistically applications beginning in 1996 with Neal’s paper on neural networks. The marriage of deterministic molecular simulations and MCMC allows fast exploration of state spaces of distributions while avoiding random walks necessary in Metropolis-Hasting algorithms.

In the subsequent sections we describe the Hamiltonian dynamics (section 3) and how to construct a MCMC method using them (section 4). Like in other Markov Chain methods HMC has its issues of tuning, which are discussed briefly in section 4. Section 5 compares HMC to Metropolis random-walk, in both low and high dimensions. In section 6 a variation of HMC is presented, where the acceptance rate of HMC is shown to be increased by looking at the “windows” of states at the beginning and the end of the trajectories.

3 Background - Hamiltonian Dynamics

Hamiltonian dynamics describes the potential energy and kinetic energy of an object, linking them through Hamilton’s Equations. The potential energy is a function of an object’s position, while the kinetic energy is the function of an object’s momentum. In HMC, the target distribution is related to the position variable, while the momentum variable is an auxiliary variable used in the deterministic updates via Hamilton’s Equations.

3.1 Hamilton’s Equations

Hamiltonian dynamics results from a certain set of differential equations. The d -dimensional vector p and q describe the full state space of $2d$ dimensions. The partial derivatives of the Hamiltonian $H(q, p)$ determine how the vectors change over time, as follows.

$$\begin{aligned}\frac{dq_i}{dt} &= \frac{\partial H}{\partial p_i} \\ \frac{dp_i}{dt} &= -\frac{\partial H}{\partial q_i}\end{aligned}\tag{3.1}$$

for $i = 1, \dots, d$. These equations map the state at any time t to a new state $t + s$, for continuous time. A discretization will be discussed in Section 3.3.

3.2 Potential and Kinetic energy

The total energy of the dynamics at point q, p of the state space is called the ‘‘Hamiltonian.’’ denoted by $H(q, p)$, and is generally considered the sum of energy generated independently by position and the momentum vector i.e.

$$H(q, p) = U(q) + K(p)\tag{3.2}$$

Here, $U(q)$ and $K(p)$ are called the potential and kinetic energy respectively. The potential energy relates to the distribution from which we want to sample, as described in Section 4.1 The kinetic energy for ease of sampling is usually defined as

$$K(p) = p^T M^{-1} p / 2\tag{3.3}$$

where M is a symmetric, positive definite matrix. As is easily seen, $K(p)$ is proportional to the multivariate Gaussian distribution with mean 0 and covariance matrix M . When M is diagonal (as is often the case), then the sampled momentum variables p are independent.

Note that $K(p)$ is a quadratic form, and we know from linear algebra that

$$\frac{dx^T A x}{dx} = x^T (A + A^T)$$

In our case, $A = M^{-1}$, which is symmetric. So $\frac{\partial K(p)}{\partial p} = p^T M^{-1}$

With these forms for H and K , Hamilton’s equation 3.1 reduces to

$$\begin{aligned}\frac{dq_i}{dt} &= [p^T M^{-1}]_i \\ \frac{dp_i}{dt} &= -\frac{\partial U}{\partial q_i}\end{aligned}\tag{3.4}$$

3.3 Discretizing Hamiltons Equation

Hamilton’s Equations describe how the system changes in continuous time - in practice, the dynamics is simulated with some small finite sample size ϵ and hence must be discretized. The discretizing methods is applicable for any form of $H(q, p)$, however we will assume the form in 3.2 as it simplifies the expressions. Also M is assumed to be diagonal with diagonal elements m_1, \dots, m_d , so that the kinetic energy has the following form.

$$K(p) = \sum_{i=1}^d \frac{p_i^2}{2m_i}\tag{3.5}$$

Which, of course, means that $\frac{\partial H(p, q)}{\partial p_i}(t) = p_i(t)/m_i$.

Euler’s method is standard method to approximate the solution to a system of differential equation. However for HMC the Leapfrog method produces better results. Thus we will skip the Euler’s method and go straight to Leapfrog method.

3.3.1 Leapfrog Method

The *leapfrog* method works as follows,

$$\begin{aligned}p_i(t + \epsilon/2) &= p_i(t) - (\epsilon/2) \frac{dp_i}{dt}(t) &= p_i(t) - (\epsilon/2) \frac{\partial U}{\partial q_i}(q(t)) \\ q_i(t + \epsilon) &= q_i(t) + (\epsilon/2) \frac{dq_i}{dt}(t + \epsilon/2) &= q_i(t) + \epsilon \frac{p_i(t + \epsilon/2)}{m_i} \\ p_i(t + \epsilon) &= p_i(t + \epsilon/2) - (\epsilon/2) \frac{dp_i}{dt}(t + \epsilon/2) &= p_i(t + \epsilon/2) - (\epsilon/2) \frac{\partial U}{\partial q_i}(q(t + \epsilon))\end{aligned}\tag{3.6}$$

The algorithm is essentially a half-step in the momentum (auxiliary) variables, a full step in the position variables using the updated momentum, and then another half-step in the momentum using the updated position. This is for one leapfrog step - typically there are many leapfrog steps per iteration. For more steps, full momentum steps are taken in between position steps and the final half-momentum step.

Because 3.6 uses mere transformations, the leapfrog method preserves volume exactly. One achieves reversibility by simply negating p due to its symmetry. These two aspects are crucial properties necessary for MCMC, but are not discussed in depth here.

4 Hamiltonian Monte Carlo

Recall that the main objective of MCMC is to simulate variables from target distributions. In our case, the variables of interest will be our momentum variables q , with target distributions uniquely determined by the potential energy function $U(q)$. Momentum variables p are introduced artificially as auxiliary random variables in order to complete the system dynamics. The following sections describe how to translate the target density into potential energy for use in the Hamiltonian, and then provide the complete HMC algorithm.

4.1 Probability and the Hamiltonian

Statistical mechanics gives us a relationship between the target distribution and potential energy function called the *canonical distribution*. Given a temperature T and energy function $E(x)$, the canonical distribution has probability density (known as the “Boltzman distribution”)

$$P(x) = \frac{1}{Z} \exp(-E(x)/T) \quad (4.1)$$

where Z is the normalizing constant.

As the Hamiltonian $H(p, q)$ defined by 3.2 is an energy function it induces a separable joint distribution over the position and the momentum as follows:

$$P(p, q) = \frac{1}{Z} \exp(-H(p, q)/T) = \frac{1}{Z} \exp(-U(q)/T) \exp(-K(p)/T) \quad (4.2)$$

It is clear from 4.2 that q and p are independent and each have canonical distribution with their own energy function. Thus the sample is generated by simulating an ergodic Markov chain that has the canonical distribution for (q, p) as its stationary distribution. We set the temperature equal to 1, unless we are using some tempering methods and define the potential energy to be

$$U(q) = -\log[P(q)] \quad (4.3)$$

where $P(q)$ is the target distribution of interest.

4.2 The Algorithm

Each iteration of the HMC algorithm has two steps. Only the momentum is changed in the first step and both the position and the momentum is changed at the second step. This explores the target density $P(q)$ defined by $U(q)$ more efficiently than using a proposal probability distribution. Starting at an initial state (q_0, p_0) , we simulate Hamiltonian dynamics for a short time using the leapfrog method. We then use the state of the position and momentum variables at the end of the simulation as our proposed states variables q^* and p^* . The proposed state is accepted using an update rule analogous to the Metropolis acceptance criterion. Specifically if the probability of the proposed state after Hamiltonian dynamics $P(q^*, p^*)$ is greater than probability of the state prior to the Hamiltonian dynamics $P(q_0, p_0)$ then the proposed state is accepted, otherwise, the proposed state is accepted randomly with probability equation to the ratio between proposed and current probabilities. If the state is rejected, the next state of the Markov chain is set as the state at $t - 1$. For a given set of initial conditions, Hamiltonian dynamics will follow contours of constant energy in phase space. Therefore we must randomly perturb the dynamics so as to explore all of $P(q)$. This is done by simply drawing a random momentum from the corresponding canonical distribution $P(p)$ before running the dynamics prior to each sampling iteration t . Combining these steps, sampling random momentum, followed by Hamiltonian dynamics and Metropolis acceptance criterion defines the HMC algorithm for drawing M samples from a target distribution:

1. set $t = 0$
2. generate an initial position state $q_0 \sim U(q)$

3. repeat until $t = M$
 - (a) set $t = t + 1$
 - (b) sample an initial momentum state $p_0 \sim K(p)$
 - (c) set $x_0 = x^{(t-1)}$
 - (d) run leapfrog starting at (p_0, q_0) for L steps and stepsize ϵ to obtain (q^*, p^*)
 - (e) calculate acceptance probability $\alpha = \min(1, \exp(U(p_0) + K(q_0) - U(q^*) - K(p^*)))$
 - (f) draw $u \sim U(0, 1)$
 - (g) if $u \leq \alpha$, accept and set $x^{(t)} = x^*$, else set $x^{(t)} = x^{(t-1)}$

HMC satisfies the “detailed balance” criterion and hence leaves the canonical distribution invariant. The algorithm is also “ergodic” and therefore explores the entire sample space.

5 Illustration of HMC and its benefits

In this section we will use several examples to demonstrate the potential of the HMC algorithm described in previous section. We will compare these to random-walk Metropolis mostly.

5.1 Trajectories for a two-dimensional problem

In this example we are sampling from a two dimensional gaussian distribution with zero mean, standard deviation one and correlation 0.95. The momentum variables are generated from an independent standard two dimensional bivariate normal. Therefore the Hamiltonian has the form

$$H(q, p) = q^T \Sigma^{-1} q / 2 + p^T p / 2, \quad \text{with} \quad \Sigma = \begin{bmatrix} 1 & 0.95 \\ 0.95 & 1 \end{bmatrix} \quad (5.1)$$

Figure(1) shows trajectories for the position coordinates and the momentum coordinates respectively along with the hamiltonian value for each of the $L = 25$ leapfrog steps with a stepsize of $\epsilon = 0.25$ that is used to generate the plots.

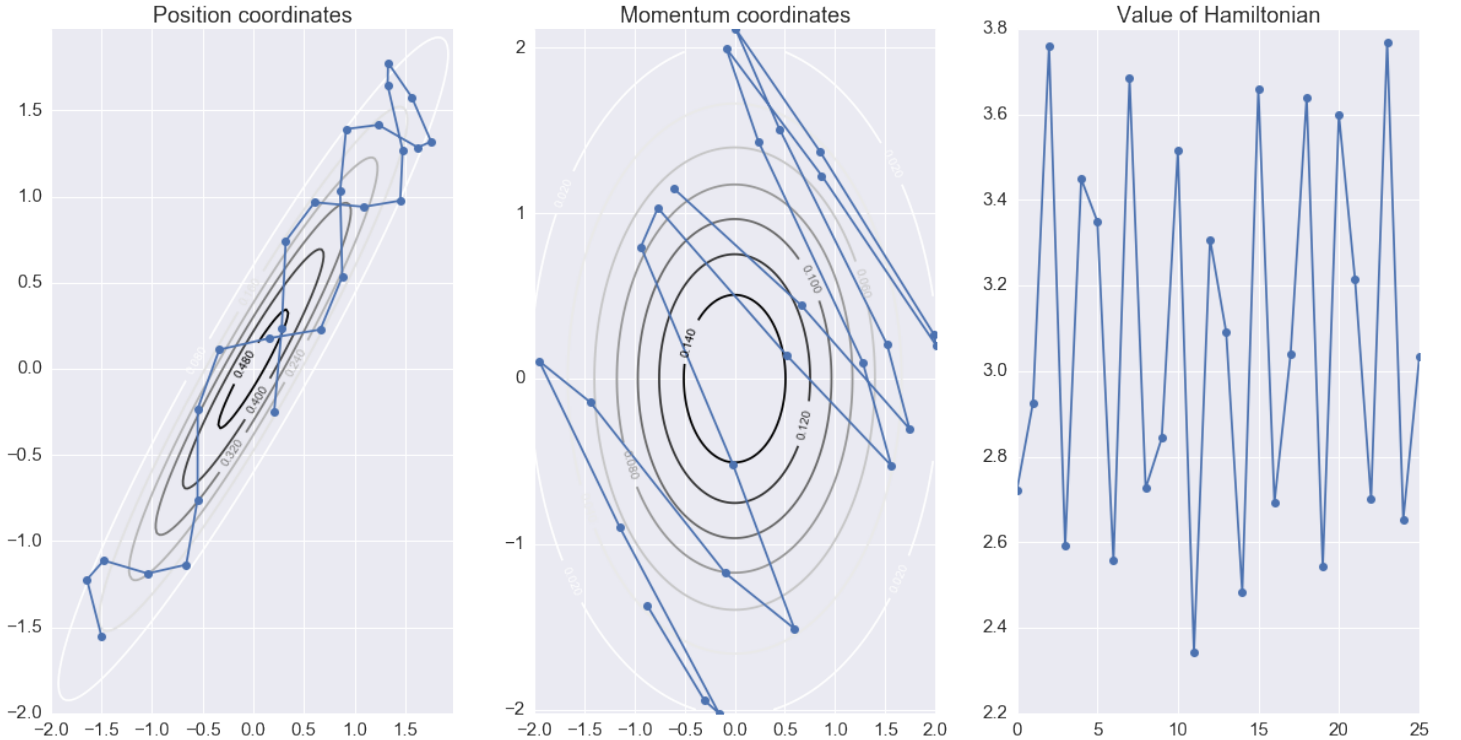


Figure 1: A trajectory for a 2D Gaussian distribution, simulated using 25 leapfrog steps with a stepsize of 0.25. The initial state for position variable q is $[-1.5, -1.5]^T$. Note how the position coordinate fully explores the space of the highly correlated distribution, while the momentum variables traverse the proposal space.

Note the differences between this and a random walk. The position variables move systematically throughout the space, reversing direction as dictated by the momentum coordinates so that the Hamiltonian oscillates appropriately. The small oscillations help the position variables fully explore the highly correlated space.

As in other methods, tuning HMC is important. Figure(2) demonstrates what can happen when the step size is catastrophically large.

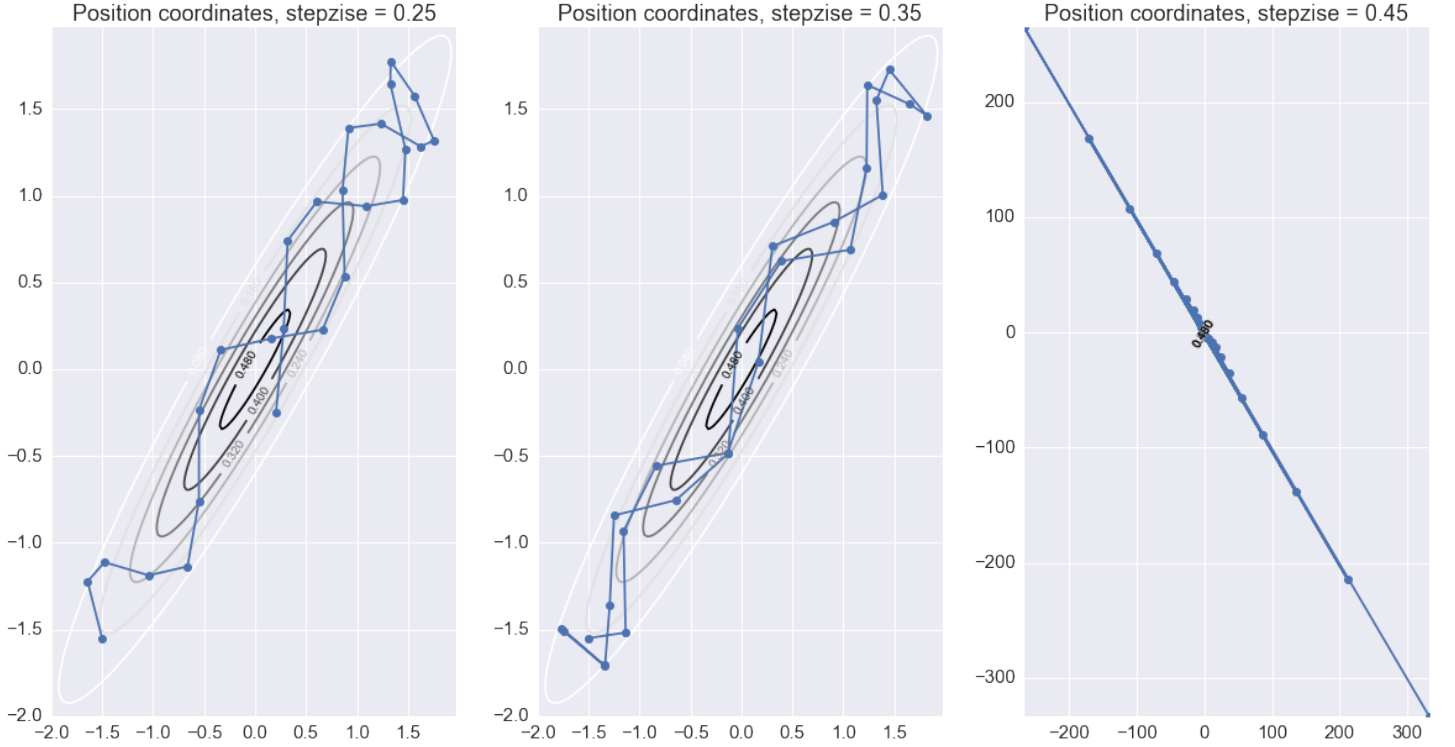


Figure 2: Illustration of Hamiltonian trajectories with increasing stepsize until it reaches critical state.

The step size is mostly controlled by the need to keep the smaller steps size under control. If larger step size were used the oscillation would be larger in the hamiltonian values. At a critical step size ($\epsilon = 0.45$), the trajectory becomes unstable and the value of Hamiltonian grows without bound (figure 2). As long as the steps size is less than that, the error in the hamiltonian stays bounded regardless of the number of leapfrog steps.

5.2 Sampling from a two-dimensional distribution

In this section we compare the trajectories between HMC and simple random-walk Metropolis method. The aim is to simulate from a bivariate Gaussian similar to the previous one but with stronger correlation of 0.98. The HMC used the same kinetic energy as before. The results of 20 HMC iteration with $L = 20$ leapfrog steps with stepsize $\epsilon = 0.18$ are shown in the right of figure(3). These values were chosen so that the trajectory length, ϵL , is sufficient to move to a distant point in the distribution without being so large that the trajectory will often back on itself. The rejection rate for this trajectories was 0.1.

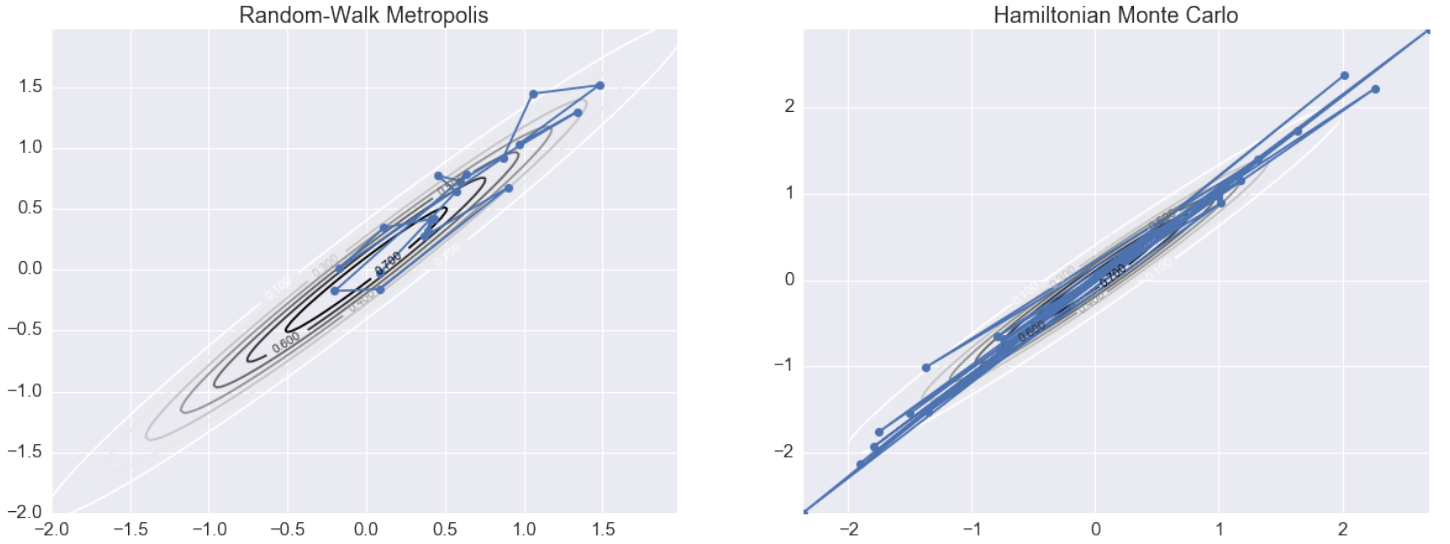


Figure 3: Twenty iterations from Metropolis (left) and HMC (right) sampling from bivariate Gaussians with .98 correlations. The Metropolis method used a thinning of 20 iterations to compare to HMC with 20 steps per iterations, but HMC still explores the space much better.

Left image of the figure(3) shows every 20th state from multiple iteration of random-walk Metropolis, with a bivariate Gaussian proposal distribution (independent with standard deviations the same as step size for HMC, $\epsilon = 0.18$). The rejection rate for the random-walk proposals is 0.38. Even with comparable steps, it is clear that HMC explores the target space much better than Metropolis, which is hampered by high correlation between the two variables.

Figure(4) illustrates the result for 200 HMC iterations and random-walk Metropolis iterations (thinned by 20 iterations). It is evident the autocorrelation between samples is much less in HMC and hence the state space is traversed more uniformly.

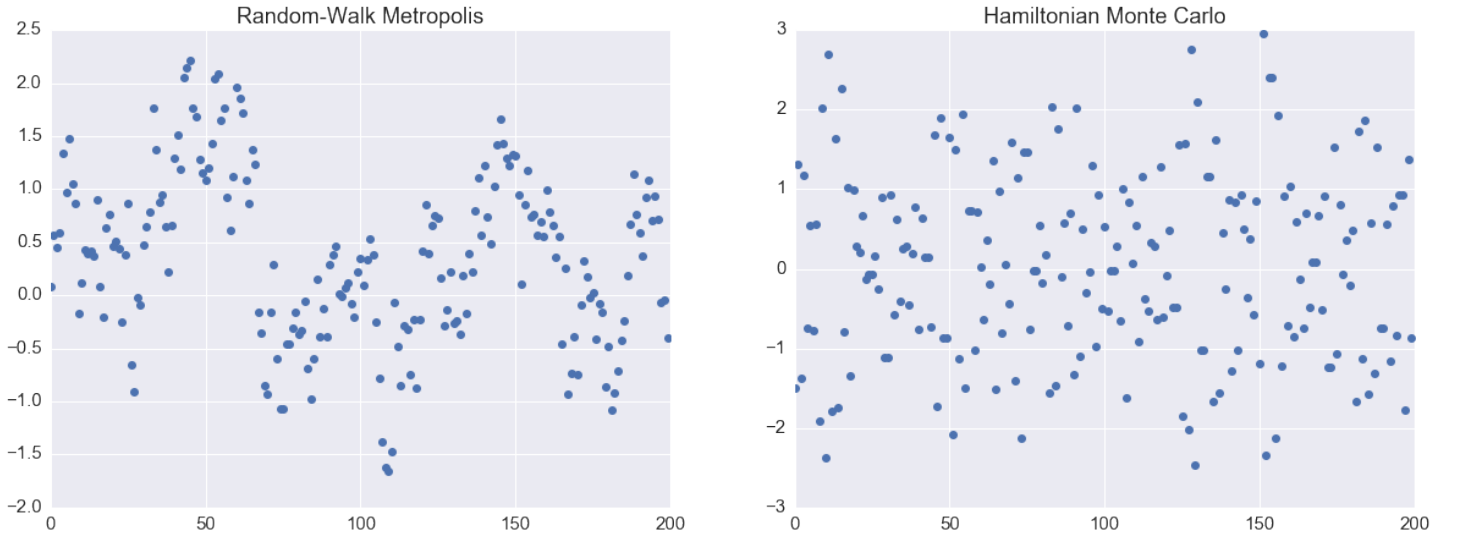


Figure 4: The first coordinate of the bivariate Gaussian sampled by the two methods. Metropolis (left) uses a thinning of 20 iterations to compare to HMC, which has 20 leapfrog steps per iteration. It is clear that there is high autocorrelation for Metropolis, even with this thinning, that is not apparent in HMC.

5.3 Efficiency in High dimension

In this section we show that HMC scales better than MH as a sampling algorithm. The target distribution to sample from is a 100 dimensional Gaussian distribution, where the variables are independent, with means of zero, and standard deviations of 0.01, 0.02, ..., 0.99, 1.00. We will use the same Kinetic energy function for the HMC as in the previous section, however it is 100 dimensional instead of 2.

For step size, we chose a number uniformly in $0.013 \pm 20\%$ - we chose a random number at each iteration in order to avoid situations where a leapfrog produces a full or half cycle of variables with standard deviations that are multiples of the step size. 0.013 was chosen to keep the step size on the same order as the smallest standard deviation of the target variables - in practice, we would not know this ahead of time, but diagnostics can be used to find the right value (selecting appropriate step sizes and number of leapfrog steps is not covered in this report). 150 leapfrog steps were used in order to allow the sampler to fully explore the space for the variables with larger standard deviations.

Figure(5) shows the sample iterations for the last dimension (and thus largest standard deviation). To compensate for the 150 leapfrog steps, the Metropolis sampler thinned by 150 iterations.

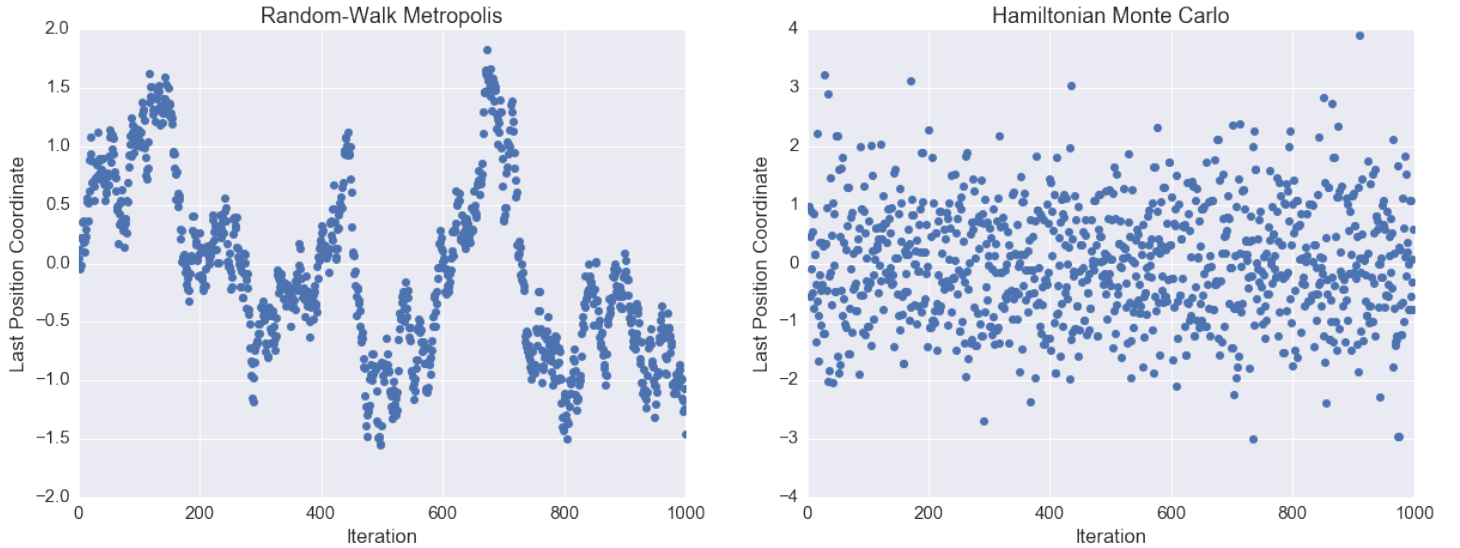


Figure 5: Iterations of sampled values for the last coordinate. The autocorrelation for random-walk Metropolis is even more pronounced.

Even when the target distributions have no correlations between variables, the high dimensionality causes the Metropolis sampler to be highly autocorrelated. Had the sampler not been thinned, the autocorrelation would have been even worse.

Perhaps more disturbing are the overall characteristics of the sampled distributions. We expect all sampled variables to have mean zero, and standard deviation to lie along the line $y=x$. Figures(6) and (7) show the results for Metropolis and HMC.

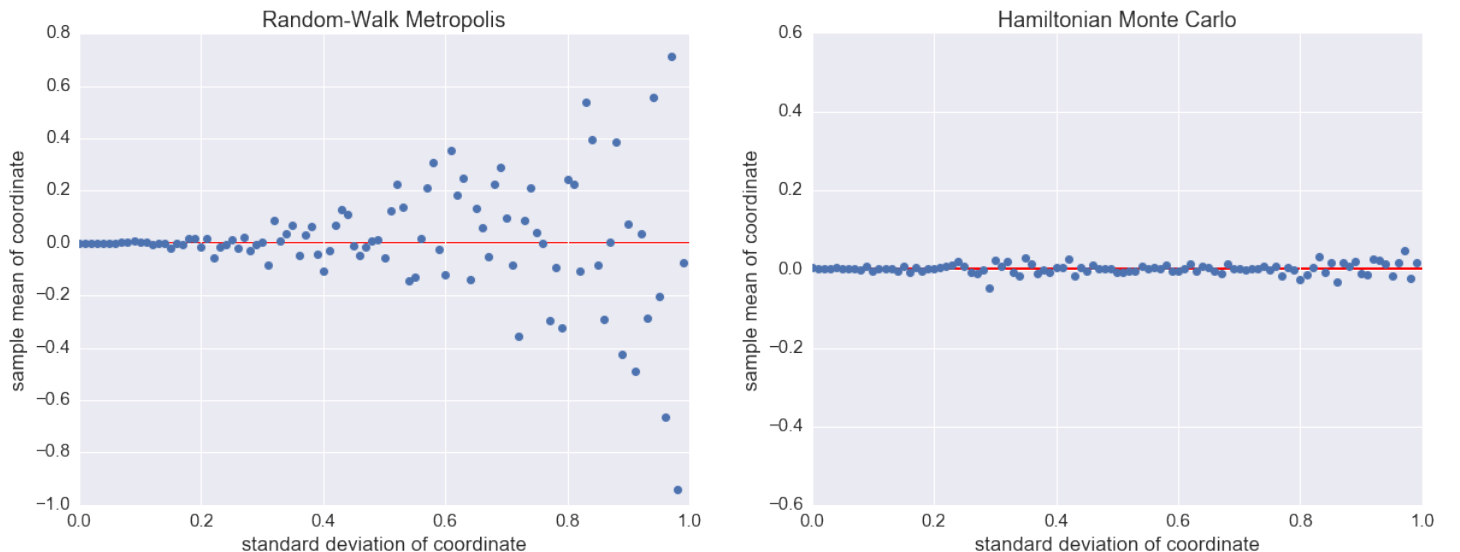


Figure 6: Sampled means for all 100 variables - we expect them all to be zero.

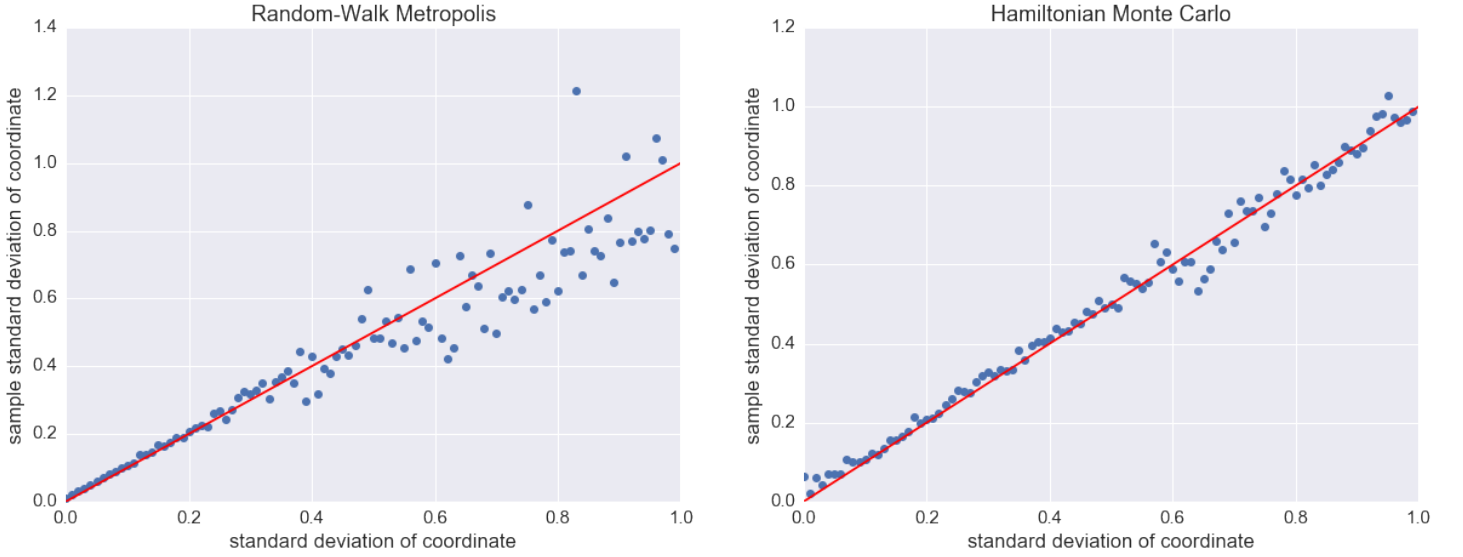


Figure 7: Sampled standard deviations for all 100 variables - we expect them to lie along the line $y=x$.

Clearly, Metropolis fails us in the higher dimensions in this case. As the dimension increases (along with the standard deviations), the means fluctuate wildly. The standard deviations of the sampled values stray from their true values, and become more erratic in the higher dimensions. HMC shows neither of these traits - all the means are close to zero, and the sampled standard deviations lie close to the line $y=x$ even in higher dimensions. Even with very simple, independent distributions, Metropolis fails in the larger dimensions where HMC succeeds just as well as in the lower dimensions.

6 Windowed States

For the final report, we plan on including a comparison to a modified version of HMC that helps obtain a higher probability of acceptance. This method, introduced by Neal in 1994, uses “windows” of states to probabilistically map the pair (q, p) to a sequence of pairs, perform similar accept-reject steps, and then probabilistically map back. Following the 1994 paper, we will show improved efficiency of the algorithm.

7 Application to Real Data

For our application to real data, we will demonstrate how the algorithm can be used to find the posterior values for parameters in simple linear regression. The model is as follows:

$$\begin{aligned} Y &\sim \mathcal{N}(X\beta, \phi^{-1}\mathcal{I}) \\ \beta &\sim \mathcal{N}(0, \tau^{-1}\mathcal{I}) \\ \phi &\sim \mathcal{G}(a, b) \end{aligned}$$

With hyperparameters $\tau, a, b = 0.1$ to keep the priors for β, ϕ reasonably flat. We are most interested in the regression parameters β - from frequentist statistics, we know that the least squares estimate $\hat{\beta} = (X^T X)^{-1} X^T Y$. If our MCMC works correctly, it should produce a posterior value centered around $\hat{\beta}$ - we can use this as a test. Before we perform HMC, however, we need the potential energy $U(\theta)$ and its gradient. In our case, $\theta = [\phi, \beta]$, with β p -dimensional. Recall that $U(\theta) = -\log(p(\theta))$ - in our case, $p(\theta) \propto \pi(\theta)L(Y | \theta)$, the posterior of $\theta | Y$. The derivation of $U(\theta)$ is below.

$$\begin{aligned} L(\theta | Y) &\propto L(Y | \phi, \beta)\pi(\beta)\pi(\phi) \\ &= \phi^{n/2} e^{-\frac{1}{2}\phi(Y-X\beta)^T(Y-X\beta)} e^{-\frac{1}{2}(\beta^T \tau \mathcal{I} \beta)} \phi^{a-1} e^{-\phi b} \\ \log(\theta | Y) &\propto \left(\frac{n}{2} + a - 1\right)\log\phi - \frac{1}{2}\phi[\beta^T X^T X \beta - 2Y^T X \beta + Y^T Y] - \frac{1}{2}\beta^T \tau \mathcal{I} \beta - \phi b \\ U(\theta) &= \left(-\frac{n}{2} - a + 1\right)\log\phi + \phi b + \frac{1}{2}\{\beta^T [\phi X^T X + \tau \mathcal{I}] \beta - 2\phi Y^T X \beta\} \end{aligned}$$

From here we can get the partial derivatives

$$\begin{aligned}\frac{\partial U}{\partial \phi} &= \frac{-n - 2a + 2}{2\phi} + b \\ \frac{\partial U}{\partial \beta} &= \beta^T [\phi X^T X + \tau \mathcal{I}] - \phi Y^T X\end{aligned}$$

References

- [1] Neal, R., M.: MCMC using Hamiltonian dynamics arXiv:1206.1901v1 5(2) (Jun 2012)
- [2] Leimkuhler, B., Reikh, S.: Simulating Hamiltonian dynamics Cambridge University Press (2004)