

# Hamiltonian Monte Carlo

## Stochastic Simulations, 2022-2023

Student Name	Student Number
Maxence Hofer	289274
Giacomo Mossinelli	359961

### 1 Introduction

Hamiltonian Monte Carlo (HMC) is a powerful tool for Bayesian computation, which offers a good computational efficiency. This method, inspired by classical mechanics, and originated for the study of lattice models of quantum field theory, soon became well known in data science. One of the key aspects of this algorithm, is given by the possibility of generating moves that can be accepted with very high probability. This acceptance probability can be brought close to 1 by implementing sufficiently accurate integration methods. In this project, we analyze this method by focusing on its theoretical framework and comparing it with other MCMC methods.

### 2 Hamiltonian conservation

Given the position vector  $\vec{q} = (q_1, \dots, q_d)$  and the momentum vector  $\vec{p} = (p_1, \dots, p_d)$ , the Hamiltonian dynamics is described by:

$$\frac{d\vec{q}_i}{dt} = \frac{\partial H}{\partial \vec{p}_i} \quad (1)$$

$$\frac{d\vec{p}_i}{dt} = -\frac{\partial H}{\partial \vec{q}_i} \quad (2)$$

for  $i = 1, \dots, d$ , where  $H(\vec{q}, \vec{p}) = U(\vec{q}) + K(\vec{p})$ , the Hamiltonian in the phase space, denotes the total energy of the system. In vector notation, it can be written also as:

$$\frac{d}{dt} \begin{bmatrix} \vec{q} \\ \vec{p} \end{bmatrix} = J^{-1} \nabla H(\vec{q}, \vec{p}) \quad (3)$$

where  $\nabla H = [\frac{\partial H}{\partial q_1}, \dots, \frac{\partial H}{\partial q_d}, \frac{\partial H}{\partial p_1}, \dots, \frac{\partial H}{\partial p_d}]^T$  and  $J = \begin{bmatrix} 0_{d \times d} & -I_{d \times d} \\ I_{d \times d} & 0_{d \times d} \end{bmatrix}$ . We can then define the trajectory

$\varphi_t : \mathbb{R}^D \rightarrow \mathbb{R}^D$ , where  $D = 2d$ , as the map which associates at the initial value  $\begin{bmatrix} q(\vec{0}) \\ p(\vec{0}) \end{bmatrix}$  the correspondent value at time  $t$  of the solution of the system of equations 1 and 2, i.e.:

$$\begin{bmatrix} q(\vec{t}) \\ p(\vec{t}) \end{bmatrix} = \varphi_t \left( \begin{bmatrix} q(\vec{0}) \\ p(\vec{0}) \end{bmatrix} \right) \quad (4)$$

Since  $J^{-1}$  is skew-symmetric, we can say that, if we consider  $(q(\vec{t}), p(\vec{t}))$  solution of 3, then:

$$\frac{d}{dt} H(q(\vec{t}), p(\vec{t})) = \nabla H(q(\vec{t}), p(\vec{t}))^T J^{-1} \nabla H(q(\vec{t}), p(\vec{t})) = 0 \quad (5)$$

This result defines the energy preservation of the Hamiltonian system, which allows us to define the following:

**Proposition 1.** *The value of the Hamiltonian function is conserved on the trajectory of the correspondent Hamiltonian system, i.e.  $H \circ \varphi_t = H \quad \forall t \in \mathbb{R}$*

Another characteristic property of an Hamiltonian system is the symplecticness, which implies that the oriented volume is conserved in the dynamics [1]. From these properties, it is possible to state that a probability measure with density of the form:

$$\mu = \frac{1}{Z} \exp(-\beta H(\vec{q}, \vec{p})) \quad (6)$$

(where  $\beta > 0$ ) is preserved over the trajectory  $\varphi_t$ . Here,  $Z$  is the normalization constant obtained with the integration  $Z = \int_{\mathbb{R}^D} \exp(-\beta H(\vec{q}, \vec{p})) dq dp < \infty$ . We can see this result as a consequence of the conservation properties since the component  $dq dp$  is preserved as a volume element and  $\exp(-\beta H(q, p))$  because it is an energy component.

### 3 Hamiltonian Monte Carlo Algorithm

Given the Hamiltonian Monte Carlo algorithm, two types of situations can be distinguished. First, let us consider the case in which the system does not require any kind of numerical discretization. In this case, the algorithm is significantly simplified, since the acceptance rate  $\alpha = \min[1, \exp(-U(\vec{q}^*) + U(\vec{q}^n) - K(\vec{p}^*) + K(\vec{p}^n))]$  is identically equal to one. More precisely, since the result given in Section 2 can be seen as a mathematical expression of the principle of conservation of energy, we can state that  $-U(\vec{q}^*) + U(\vec{q}^n) - K(\vec{p}^*) + K(\vec{p}^n) = 0$  and hence the result follows directly.

This fact changes when one considers the case where the system is subject to a numerical approximation. In fact, since in most cases the exact solution  $\varphi_t$  is not available, it is necessary to introduce the approximation  $\Psi$  but, since it is difficult to conserve energy and the volume in phase state,  $H$  can be only approximately invariant [4]. As a consequence, we are introducing an approximation error that cannot be ignored and the acceptance rate becomes:

$$\alpha(\vec{q}, \vec{p}) = \min[1, \exp(-\Delta H(\vec{q}, \vec{p})) |det \Psi'(\vec{q}, \vec{p})|] \quad (7)$$

with  $\Delta H = H(\Psi(\vec{q}, \vec{p})) - H(\vec{q}, \vec{p})$ . There are many methods that can be used to describe the evolution of the Hamiltonian system, such as the simple Euler method, but the Verlet method is one of the most widely used in this context which, in our case, is:

$$p_i(t + \frac{\epsilon}{2}) = p_i(t) - \frac{\epsilon}{2} \frac{\partial U(q(t))}{\partial q_i} \quad (8)$$

$$q_i(t + \epsilon) = q_i(t) + \epsilon \frac{p_i(t + \frac{\epsilon}{2})}{m_i} \quad (9)$$

$$p_i(t + \epsilon) = p_i(t + \frac{\epsilon}{2}) - \frac{\epsilon}{2} \frac{\partial U(q(t + \epsilon))}{\partial q_i} \quad (10)$$

This method presents many advantages. The global error of the method has order 2 (unlike the aforementioned Euler method, of order 1) [1]. Furthermore, being a symplectic integrator, it is volume preserving and, consequently, it allows a simpler expression to define  $\alpha$ , not involving the determinant of the trajectory approximation ( $\alpha(\vec{q}, \vec{p}) = \min[1, \exp(-\Delta H(\vec{q}, \vec{p}))]$ ). However, the transition from the exact form to the numerical form implies the loss of the principle of conservation of energy and, consequently, does not allow one to state that  $\alpha$  is identically equal to 1.

## 4 Exact time integration

In this section, in order not to have excessively heavy writing, we have avoided putting the arrow to indicate that  $q$  and  $p$  are a vector. Given the potential energy  $U(q) = -\log \pi(q)$  and the kinetic energy  $K(p) = \sum_{i=1}^d \frac{p_i^2}{2m_i}$ , we consider the Gibbs-distribution as:

$$G(q, p) = \frac{1}{\tilde{Z}} \exp(U(q)) \frac{1}{\prod_{i=1}^d \sqrt{2\pi m_i}} \exp(-K(p)) \quad (11)$$

where  $q$  is the variable of interest, while  $p$  is introduced only to allow Hamiltonian dynamics to operate. In fact, Hamiltonian Monte Carlo uses the Hamiltonian dynamical system to construct the Markov Chain Monte Carlo algorithm on a given invariant density  $\pi(q)$ .

If we consider a case where the time integration is being carried-out exactly, it is possible to show that equation 11 is invariant [1]. To do so, it is necessary to introduce the reversibility with respect to the momentum flip involution of the Hamiltonian system, which is related to the fact that  $H(q, p)$  is even w.r.t.  $p$ . This property, which can also be seen as the time-reversibility of mechanics, states that, if  $(q^0, p^0)$  is the initial state of the system and  $(q(T), p(T))$  is the final state after a timespan  $T$ , then  $(q(T), -p(T))$  evolves in the same amount of time to the state  $(q^0, -p^0)$ .

Following the Hamiltonian Monte Carlo algorithm, we can see that  $p^0 \sim \mathcal{N}(0, M)$ , which obviously preserves Gibbs distribution. Then, we let the system evolve deterministically for a fixed time  $T$ , according to Hamilton's system (1,2). This evolution defines a mapping  $(q^0, p^0) \rightarrow (q^*, p^*) = (q(T), p(T))$  which obviously implies that  $\alpha = 1$ . Since this process is deterministic, we can also say that  $g(q^*, p^*, q^0, p^0) = P[(q(T), p(T)) = (q^*, p^*) | (q(0), p(0)) = (q^0, p^0)] = \delta((q(T), p(T)) - (q^*, p^*))$ . The transition probability for  $q$  can then be defined as:

$$P[q(T) = q^* | q(0) = q^0] = \int P[p(0) = p^0] g(q^*, p^*, q^0, p^0) \alpha(q^*, p^*)$$

and we would like to prove that it is in detailed balance [2].

Since

$$\exp(-H(q, p)) \min(1, \exp(-\Delta H)) = \exp(-H(q^*, p^*)) \min(\exp(\Delta H), 1)$$

and since the Gibbs distribution  $G(q, p)$  has the same form as equation 6, and so depends on  $H(q, p)$ , which is reversible, we obtain:

$$G(q^0, p^0) \alpha(q^*, p^*) = G(q^*, p^*) \alpha'(q^0, p^0) = G(q^*, -p^*) \alpha'(q^0, -p^0)$$

where  $\alpha'(q^0, -p^0)$  denotes the acceptance in the opposite direction. Finally, using the fact that  $g(q^*, p^*, q^0, p^0) = g(q^0, -p^0, q^*, -p^*)$  because of reversibility, then we have

$$\int G(q^0, p^0) g(q^*, p^*, q^0, p^0) \alpha(q^*, p^*) = \int G(q^*, -p^*) g(q^0, -p^0, q^*, -p^*) \alpha'(q^0, -p^0) \quad (12)$$

which proves the detailed balanceness. In the case where both the energy and the volume are conserved as in paragraph 2, then  $\alpha = 1$  always. This proves that the distribution will be correct for any configuration of  $q$  and that the Gibbs distribution is conserved.

The same is not true when we use the Verlet method, since in this case we introduce an error in the energy conservation.

## 5 Implementation of the algorithm

Markov Chain Monte Carlo algorithm are used in order to sample from a unnormalized distribution  $\pi$ . The idea is to create an ergodic Markov Chain with invariant distribution  $\pi$  [5]. In

the following part we will compare the performance of two different algorithms to sample from the unnormalized distribution:

$$f_1(q_1, q_2) = e^{-\alpha(q_1^2 + q_2^2 - 1/4)^2} \quad (13)$$

More precisely, in section 6 we will implement a Hamiltonian Monte Carlo method (HMC) using a Verlet's scheme to simulate the next element of the chain and a Random Walk Metropolis-Hastings (RWMH) algorithm.

**Random Walk Metropolis-Hastings** is an Markov chain Monte Carlo algorithm using a Normal distribution to simulate the next point  $\vec{q}^* \sim \mathcal{N}(\vec{q}_n, \sigma^2 I_2)$ , with  $\vec{q}_n$  the last element of the chain and  $\sigma$  a parameter to adjust. This point  $\vec{q}^*$  will then be accepted with a probability  $\beta$  defined as

$$\beta = \min \left\{ \frac{f(\vec{q}^*)}{f(\vec{q}_n)}, 1 \right\} \quad (14)$$

If the next element is accepted,  $q_{n+1} = \vec{q}^*$ , otherwise  $q_{n+1} = \vec{q}_n$ .

In order to analyse the Markov chains produced by both algorithms, we need to define two usefull functions.

First, the autocovariance of a Markov chain is a measure of the degree to which the current state of the chain is correlated with its past states. It is defined as:

$$c(k) = \text{Cov}(X_j, X_{j+k}) \quad (15)$$

where  $X_t$  is the state of the Markov chain at time  $t$ , and  $k$  is the lag at which the autocovariance is being computed.

This function can be used to assess the persistence of the chain over time. If the autocovariance decays quickly as the lag increases, then the chain is said to have low persistence, which means that the current state is largely independent of the past states of the chain. On the other hand, if it remains high for large lags, then the chain is said to have high persistence, which means that the current state is strongly influenced by the past states of the chain.

The autocovariance function is often used in the study of Markov chain, and it is closely related to the autocorrelation function, which is another measure of the correlation between the current state of the chain and its past states. The autocorrelation function is defined as:

$$\rho(k) = \frac{\text{Cov}(X_j, X_{j+k})}{\text{Var}(X_j)} \quad (16)$$

where  $\text{Var}(X_t)$  is the variance of the state at time  $t$ .

The effective sample size (ESS) of a Markov chain represents the size of an equivalent independent sample that would lead to the same variance of the estimator. This is defined as:

$$\text{ESS} = N \frac{c(0)}{\sigma_{mcmc}^2} \quad (17)$$

Where  $N$  is the length of the chain,  $c(0)$  the variance of  $X_j$  and  $\sigma_{mcmc}^2$  is called time-average variance constant and is given by:

$$\sigma_{mcmc}^2 = \text{Var}(X_j) + 2 \sum_{l=1}^{\infty} \text{Cov}(X_j, X_{j+l}) < \infty \quad (18)$$

In our case, we have calculated an estimator of the covariance as:

$$\hat{c}(k) = \frac{1}{N - k - 1} \sum_{j=B+1}^{N+B-k} (X_j - \hat{\mu}_{N,B}^{mcmc})(X_{j+k} - \hat{\mu}_{N,B}^{mcmc}) \quad (19)$$

with  $B$  a burn-in lag during which the chain reaches a stationary state, and  $\mu_{N,B}^{mcmc}$  the average of the chain after the burn-in lag. the correlation as

$$\hat{\rho}(k) = \frac{\hat{c}(k)}{\hat{c}(0)} \quad (20)$$

and the time-average variance constant

$$\hat{\sigma}_M^2 = \hat{c}(0) + 2 \sum_{k=1}^M \hat{c}(k) \quad (21)$$

where we only use the first terms of the sum, because of the instability of the last terms.  $M$  is defined as  $M = 2\min\{k : \hat{c}(2k) + \hat{c}(2k+1) < 0\}$ .

In this experiment parameters of both algorithm are the following:

1. **HMC:**

- $\epsilon$ : the step size in the Verlet's scheme
- $T$ : the final time of integration in the Verlet's scheme
- $\vec{m}$  the vector of the mass used in the scheme

2. **RWMH:**

- $\sigma^2$  the variance of the Normal distribution used to compute  $\vec{q}^*$

3. **Both algorithms:**

- $\alpha$ : the constant in the distribution, that will be studied for  $\alpha = 10$  and  $\alpha = 1000$
- $N$ : the length of the chains
- $n$ : the number of chains simulated
- $\vec{q}_0$ : the distribution of the initial state of the chain, in this exercise, we will use  $\vec{q}_0$  a Dirac distribution centered in  $[0, 0]$

## 6 Results

In this section we present the results obtained with the implemented algorithms. First of all, we used both RWMH and HMC algorithms to simulate  $n=1000$  chains. Figure 1 represents the final  $q$  of the chains, with initial parameters  $\alpha=10$ ,  $\vec{q}_0 = \delta([0, 0])$ ,  $m_1 = 1$ ,  $m_2 = 1$ ,  $N = 100$ ,  $n = 1000$ ,  $\epsilon = 0.01$ ,  $T = 0.1$ ,  $\sigma^2 = 0.1$ . Both algorithms use a Dirac distribution centered in

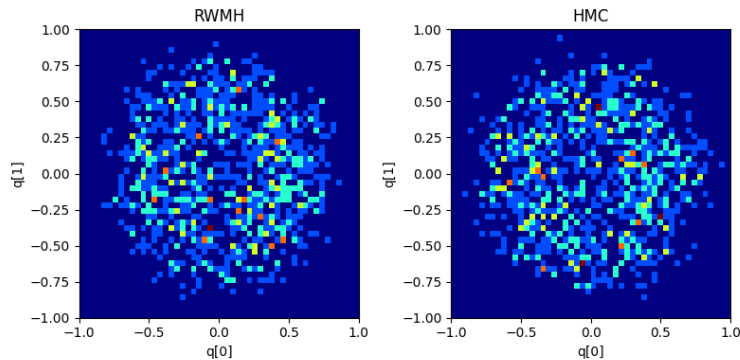


Figure 1: Final  $\vec{q}$

$[0, 0]$  to produce the Markov chains. Once the stationary state is reached, the final distribution is invariant under the effect of both algorithms. Therefore, it is possible to plot the sum of every element of the Markov chain after a certain burn-in lag  $B$ , when the chain has reached its stationary state. Figure 2 represents the density after burn-in lag  $B = 20$ , using the previously mentioned parameters. The left figure represents the RWMH algorithm, the central figure the

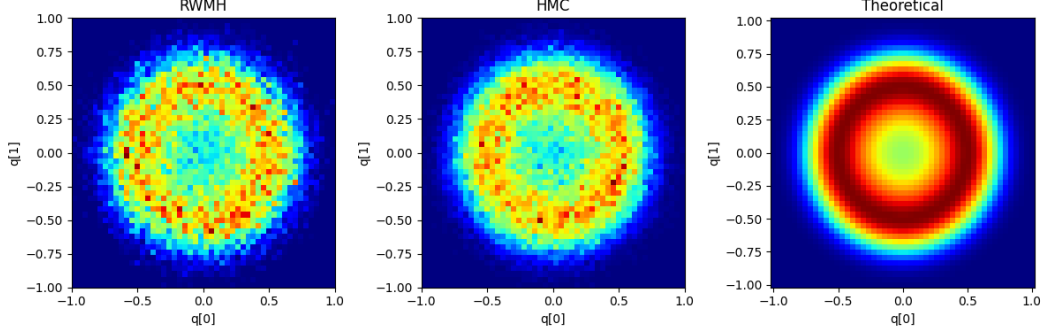


Figure 2: Sum of  $\vec{q}$  starting at  $B=20$ ,  $\alpha = 10$

HMC method and the right one the distribution  $f$ . It is possible to observe that both methods seem to model correctly the distribution. For what concerns RWMH, this result was expected also from a theoretical point of view. Since the distribution is log-concave in tails, the Markov chain generated by the algorithm should be geometrically ergodic [5]. The ratio of acceptance is for RWMH  $r_{RWMH} = 0.632$  and for HMC  $r_{HMC} = 0.9997$ . This result for HMC is in line with what was explained in section 3. Since we are exploiting a symplectic integrator, the problem is approximated with an high accuracy, even if it was not possible to reach a perfect acceptance rate. Figure 3 represents the same process, but with  $\alpha = 10^3$ ,  $B = 80$  and  $N = 200$ . In this case, the two methods seem to converge to the distribution  $f$ .

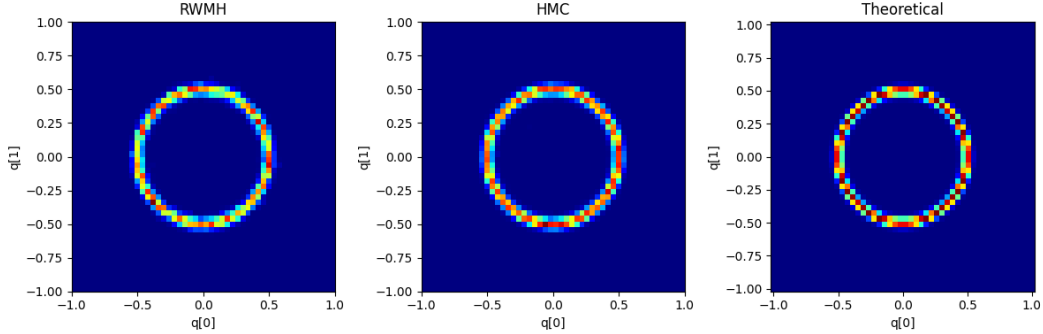


Figure 3: Sum of  $\vec{q}$  starting at  $B=80$ ,  $\alpha = 10^3$

We can observe the consequences of the modification of some of the parameters on the HMC method. Fig. 4, for example, represents the effect of different  $\vec{m}$ ,  $T$ , and  $\epsilon$ . In this case, we adopted:  $\alpha = 10$ ,  $\vec{q}_0 = \delta([0, 0])$ ,  $N = 100$ ,  $n = 1000$ . It is possible to observe that an increase of  $m_i$ ,  $\epsilon$  or a decrease of  $T$  imply that the simulation does not follow the distribution  $f$ . This is probably due to the fact that, with these parameters, the chain does not reach a stable state. It is then possible to study different properties on these Markov chains. Figures 5 and 6 represent the average and variance on the 1000 chains, as a function of the state of the chains, respectively for  $\alpha = 10$  and  $\alpha = 10^3$ . We can observe that, as the initial distribution is a Dirac distribution centered in  $[0, 0]$ , the initial value of average and variance are both zeros. Then, the average oscillates around zero while the variance stabilizes around 0.15 for  $\alpha = 10$  and 0.12 when  $\alpha = 10^3$ . Another meaningful information we can deduce from this graphs is the instant when the stability of the variance occurs. More precisely, we can define the beginning of this

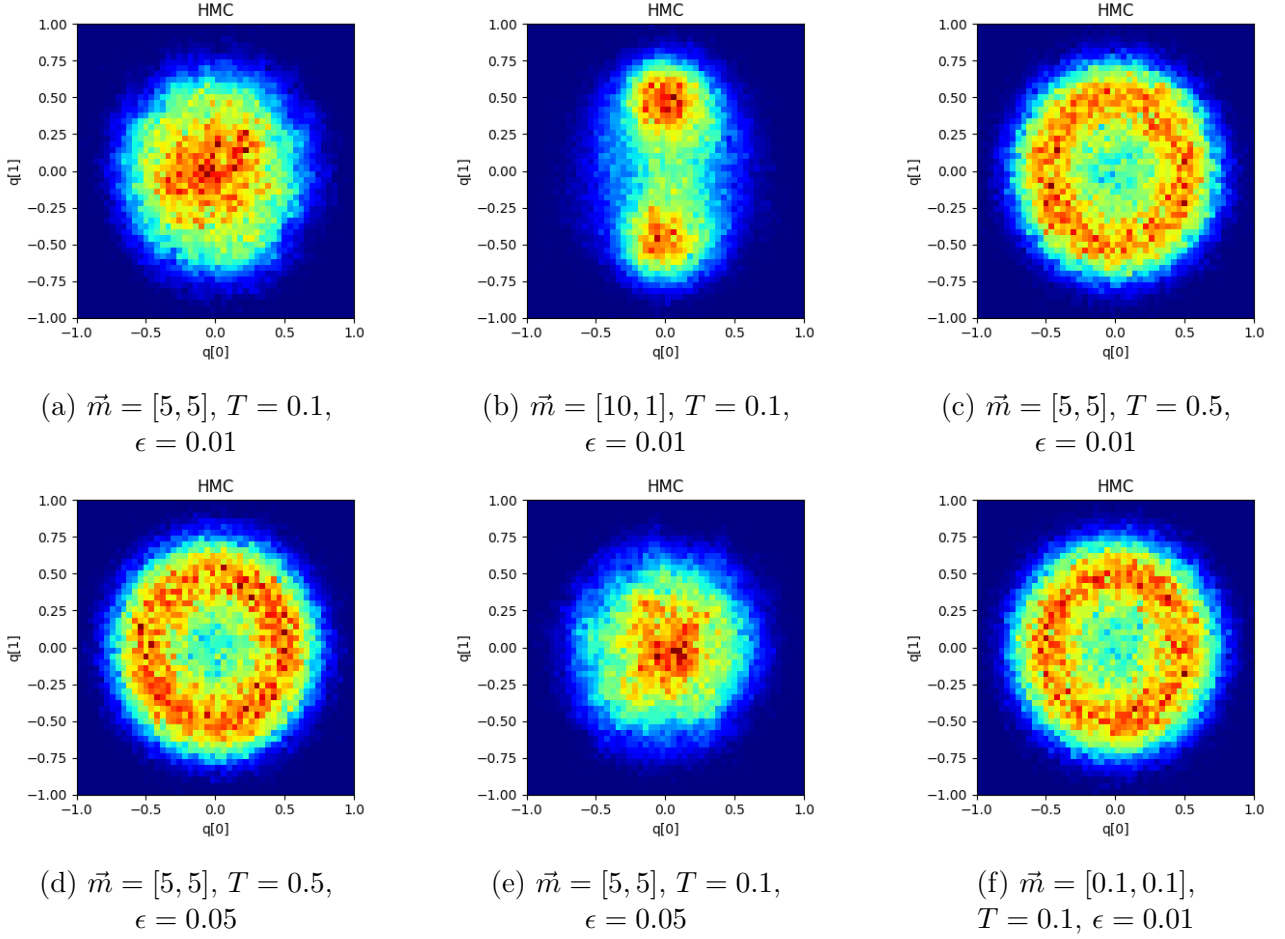


Figure 4: Effect of different parameters for HMC

trend before 10 elements of the Markov chain for both methods when  $\alpha = 10$ , and around 80 for RWMH and HMC when  $\alpha = 10^3$ . These values allow us to define the burn-in lag, which will be used to calculate the autocovariance of the chain.

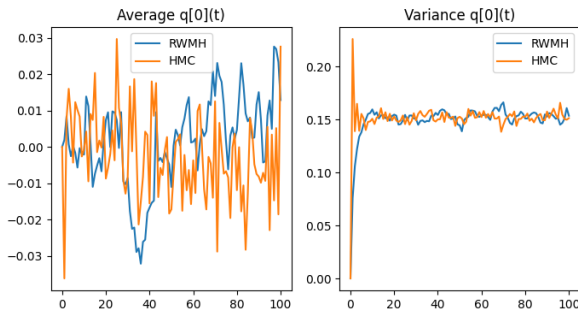


Figure 5: Average and variance of  $\vec{q}$ ,  $\alpha = 10$

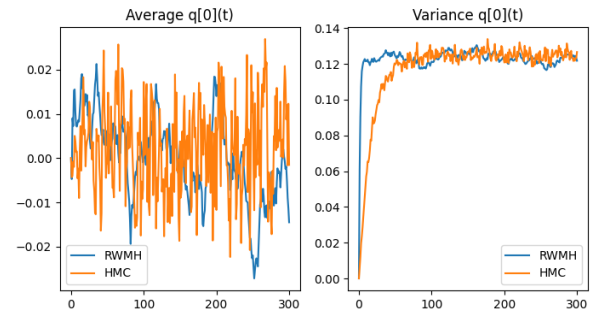


Figure 6: Average and variance of  $\vec{q}$ ,  $\alpha = 10^3$

Other key features for a detailed analysis, as more deeply explained in section 5, are the autocovariance and autocorrelation. In this case, the average correlation on  $n$  chains in the HMC chain is going quickly to zero, which means that each state is largely independent of the past states of the chain. For what concern RWMH, on the contrary, the autocorrelation is not following the exact same path. This may mean zero there is a higher relation between the states. Therefore, it is possible to conclude that the sampling from HMC has more potential, as it provides more independent samples than RWMH.

These results are represented in figures 7 and 8 while in figures 9 and 10 the averages on the

n chains are reported. Since there is a symmetry between  $q[0]$  and  $q[1]$ , the graphs are only shown for  $q[0]$ . We also observe that the autocorrelation and autocovariance for one chain are not stable. This is probably due to the calculation of the covariance, since the elements  $\hat{c}(k)$  are samples average of very few terms. Parameters used are:  $\vec{q}_0 = \delta([0, 0])$ ,  $\vec{m} = [0.1, 0.1]$ ,  $n = 1000$ ,  $\epsilon = 0.01$ ,  $T = 0.2$ ,  $\sigma^2 = 0.1$ ,  $N = 100$ ,  $B = 20$  for  $\alpha = 10$  and  $N = 500$ ,  $B = 100$  for  $\alpha = 10^3$ .

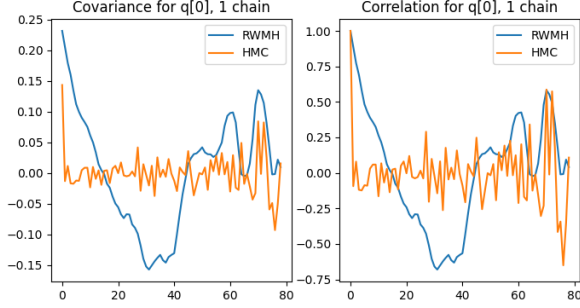


Figure 7: Covariance and correlation of 1 chain for  $\vec{q}$ ,  $\alpha = 10$

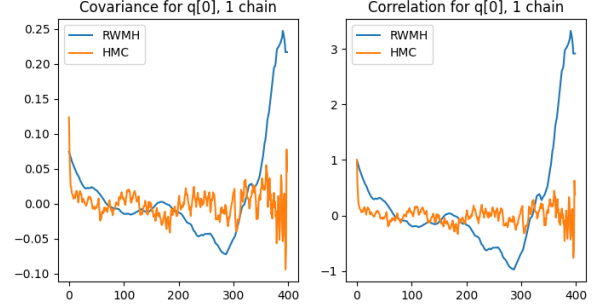


Figure 8: Covariance and correlation of 1 chain for  $\vec{q}$ ,  $\alpha = 10^3$

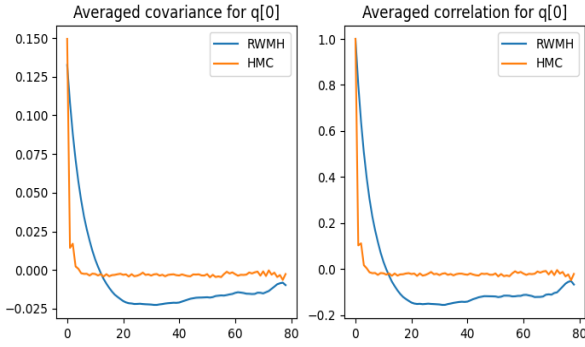


Figure 9: Averaged covariance and correlation for  $\vec{q}$ ,  $\alpha = 10$

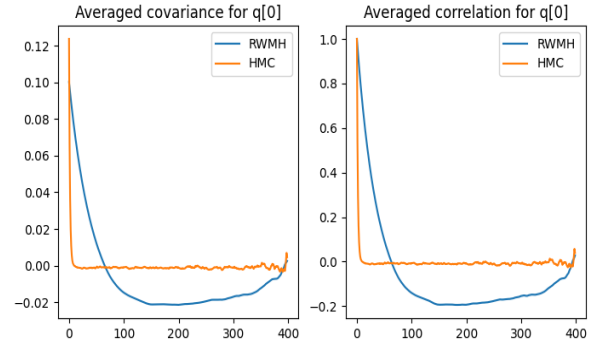


Figure 10: Averaged covariance and correlation for  $\vec{q}$ ,  $\alpha = 10^3$

The previously mentioned remark is also shown in figures 11 and 12, that represent the distribution of the effective sample size (ESS), as described in section 5 obtained using the same parameters as before. We observe, as we could expect from the autocorrelation graphs, that the effective sample size is higher for HMC than for RWMH.

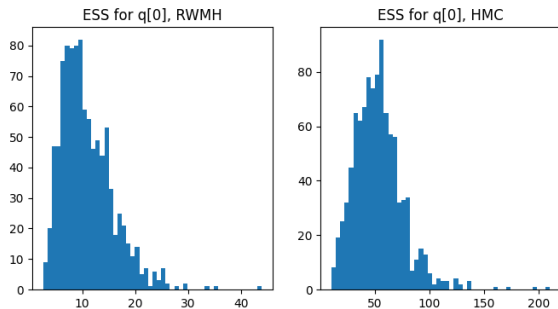


Figure 11: Distribution of effective sample size,  $\alpha = 10$

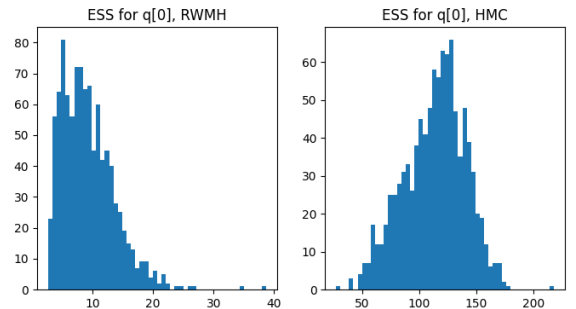


Figure 12: Distribution of effective sample size,  $\alpha = 10^3$



Finally, one last useful thing to analyse is the effective sample size as a function of the number of evaluations of  $f$  and  $\partial f/\partial q_i$ . The results are shown in figures 13 and 14. Parameters are given by:  $\vec{q}_0 = \delta([0, 0])$ ,  $\vec{m} = [0.1, 0.1]$ ,  $n = 50$ ,  $\epsilon = 0.01$ ,  $T = 0.2$ ,  $\sigma^2 = 0.1$ ,  $N = 50 - 500$ ,  $B = 20$  for  $\alpha = 10$  and  $N = 100 - 500$ ,  $B = 80$  for  $\alpha = 10^3$ . The length of the chain  $N$  has been varied between the two reported values to get different number of evaluations of  $f$  and  $\partial f/\partial q_i$ . The number of evaluations  $n_{ev}$  for HMC is equal to  $n_{ev} = 2 + 4\frac{T}{\epsilon}N$  while, for RWMH, is defined as  $n_{ev} = 2N$ . One can see that the effective sample size is linearly increasing as a function of  $n_{ev}$  of  $f$  and  $\partial f/\partial q_i$ . For both algorithms, the effective sample size is higher for  $\alpha = 10$  than for  $\alpha = 10^3$ . The number of evaluations of  $f$  and  $\partial f/\partial q_i$  is higher for HMC than for RWMH, due to the fact that the Verlet's scheme needs each time  $O(4 * \frac{T}{\epsilon})$  evaluations when the RWMH only needs two computation of  $f$  in order to calculate the coefficient of acceptance. But as the Verlet's scheme does not require a high computational time, it would not make sense to compare the performance of both algorithms for the same number of evaluations of  $f$  and  $\partial f/\partial q_i$ .

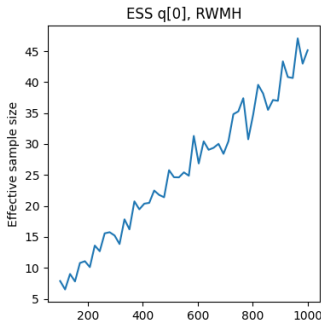


Figure 13: ESS as a function of the number of evaluation of  $f$  and  $\partial f/\partial q_i$ ,  $\alpha = 10$

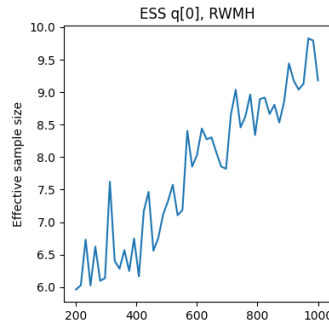
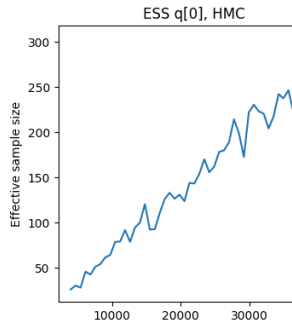
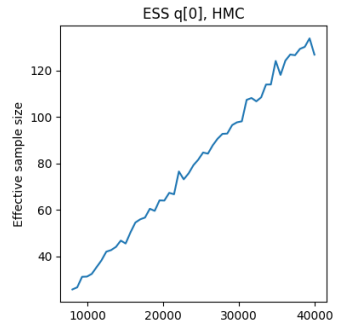


Figure 14: ESS as a function of the number of evaluation of  $f$  and  $\partial f/\partial q_i$ ,  $\alpha = 10^3$



## 7 Application

In this section, we applied our Hamiltonian Monte-Carlo to sample from a distribution reported in [3]. The dataset, contains data of 189 births, collected at Baystate Medical Center in Springfield, Massachusetts. More precisely, there are information about the weights of newborns and indications about characteristics of the mother which can be correlated to this. It was asked to perform a logistic regression with the HMC sampler on this dataset.

First of all, the proposed dataset has been slightly modified to meet the requests. For example, a column for the intercept has been added and the initial covariate "RACE" (which had value 1 for white, 2 for black and 3 for other) has been splitted in 2 categorical binary features, "race2black" and "race2other". Moreover, the column "LOW" as been used as binary response of our model. The features of the dataset are:

- **intersept**: column of ones added to the original dataset
- **age**: integer values which denotes the age of the mother
- **lwt**: mother's weight in pounds at last mestrual period
- **race2black**: categorical variable: 0 if white, 1 if black
- **race2other**: categorical variable: 0 if white or black, 1 if other
- **smoke**: categorical variable: 1 if smoking during pregnancy

- **ptd**: categorical variable: 1 if premature birth
- **ht**: categorical variable: 1 if hypertension
- **ui**: categorical variable: 1 if uterine irritability
- **ftv21**: categorical variable: 1 if one physician visit during the first trimester
- **ftv22+**: categorical variable: 1 if at least two physician visits during the first trimester

In particular, the two continuous variables (namely, age and lwt) have been normalized to assume values more comparable to the other features. The binary response which divides babies into underweight and normal-weight categories is connected to the other elements of the dataset through a logistic regression model  $P(y_i = 1|x_i^T q) = [1 + \exp(-x_i^T q)]^{-1}$  with  $i = 1, \dots, n$  where  $n$  is the number of data. This problem can be modelled implementing the HMC algorithm. Even if this method has a physic background, it can be easily generalized to be suitable in a lot of different contexts. A description of this fact is, for example, reported in [6]. In this case, we have 11 covariates instead of the 2 we have considered in the previous analysis. Moreover, we have to specify a log-posterior for the model, which in this case is:

$$\log f(q|y, X, \sigma_q^2) = q^T X^T (y - \mathbf{1}_n) - \mathbf{1}_n^T [\log(1 + \exp(-x_i^T q))]_{n+1} - \frac{q^T q}{2\sigma_q^2} \quad (22)$$

with gaussian prior  $\mathcal{N}(0, \sigma^2 I_{p+1})$  where  $p+1$  is the number of the covariates. In the Hamiltonian Monte Carlo algorithm, instead of the potential energy, we have then

$$H(q, p) = \log(q|y, X, \sigma_q^2) + \frac{1}{2} p^T M^{-1} p \quad (23)$$

The Verlet method is obviously modified too, and it becomes

$$p_i(t + \frac{\epsilon}{2}) = p_i(t) + \frac{\epsilon}{2} \nabla_q \log f(q(t)|y, X, \sigma_q^2) \quad (24)$$

$$q_i(t + \epsilon) = q_i(t) + \epsilon \frac{p_i(t + \frac{\epsilon}{2})}{m_i} \quad (25)$$

$$p_i(t + \epsilon) = p_i(t + \frac{\epsilon}{2}) + \frac{\epsilon}{2} \nabla_q \log f(q(t + \epsilon)|y, X, \sigma_q^2) \quad (26)$$

The parameters used to model the problem are the following:  $\vec{q}_0 = \delta(\vec{0})$ ,  $\vec{m} = [\vec{1}]$ ,  $n = 10$ ,  $\epsilon = 0.01$ ,  $T = 0.4$ ,  $\sigma^2 = 1000$ ,  $N = 400 - 6000$ ,  $B = 200$ . These parameters were chosen because the results provided using them were satisfying. In order to study the convergence of the HMC method, we have analysed the effective sample size according to the length of the chain. Since the ESS is a measure of the number of independent samples that are being drawn from the target distribution, an increasing value of this parameter would imply a convergence of the method. In fact this phenomenon means that the samples are becoming less correlated and more representative of the target distribution. Fig. 15 represents the ESS for each feature. We observe that the effective sample size is globally increasing as the length of the Markov chain grows with no exception, therefore implying a convergence of the HMC method. We present then, in figure 17 the histograms of the simulated posteriors obtained using the previously mentioned parameters.

Finally, we compare the performance of Hamiltonian Monte Carlo with a one variable at a time Metropolis-Hastings algorithm (1MH) to compute the posterior expectation of  $q$ . One Variable at a time Metropolis-Hastings is similar to RWMH presented in sec. 5 but modifying only one feature at each step, using variance of the normal distribution  $\sigma^2 = 0.5$ . Figure 16 represents the average on  $n = 50$  chains of  $q(t)$  for 4 different features. The parameters were:  $\vec{q}_0 = \delta(\vec{0})$ ,

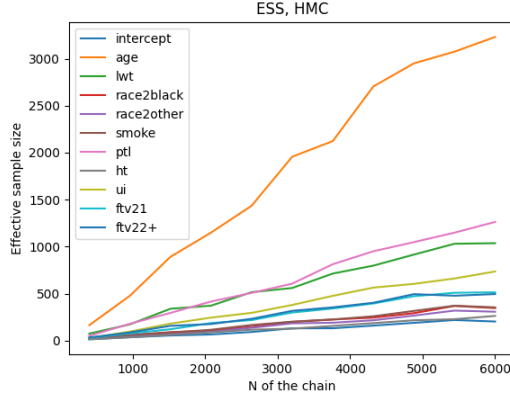


Figure 15: Effective sample size for every feature, as a function of  $N$

$\vec{m} = [\vec{1}]$ ,  $n = 10$ ,  $\epsilon = 0.01$ ,  $T = 0.4$ ,  $\sigma^2 = 1000$ ,  $N = 400$ . We observe that the chains produced by HMC algorithm stabilize earlier than 1MH, for any feature and regardless of the value of the variance of the 1MH algorithm. This was expected because the HMC method calculate the next point using the gradient of the log-posterior, while the proposition of next element of the chain using MH method is sampled according to a normal distribution centered in the previous point. Moreover, it is important to observe that the 1MH algorithm validates the results obtained with HMC since both methods seems to reach convergence in correspondence of the same values.

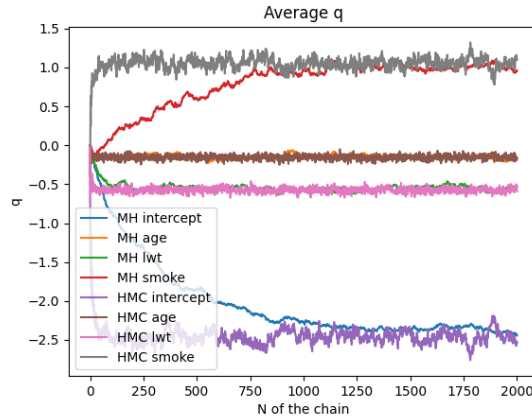


Figure 16: Average on  $n$  chains of different features

## 8 Conclusion

In this project, we presented an implementation of a Hamiltonian Monte Carlo method. After describing some theoretical results in sections 2, 3, 4 and methodology to implement the algorithm in section 5, we applied our algorithm to sample from an unnormalized distribution. Comparing HMC with Random Walk Metropolis-Hastings, we have shown that the HMC performed better to sample from the given distribution. Then, we applied the HMC to perform a logistic regression task in section 7, and we observed that the method was convergent. Finally, we compared the HMC for this problem with a One Variable at a time Metropolis-Hastings and we also observed that the HMC performed better than the other method.

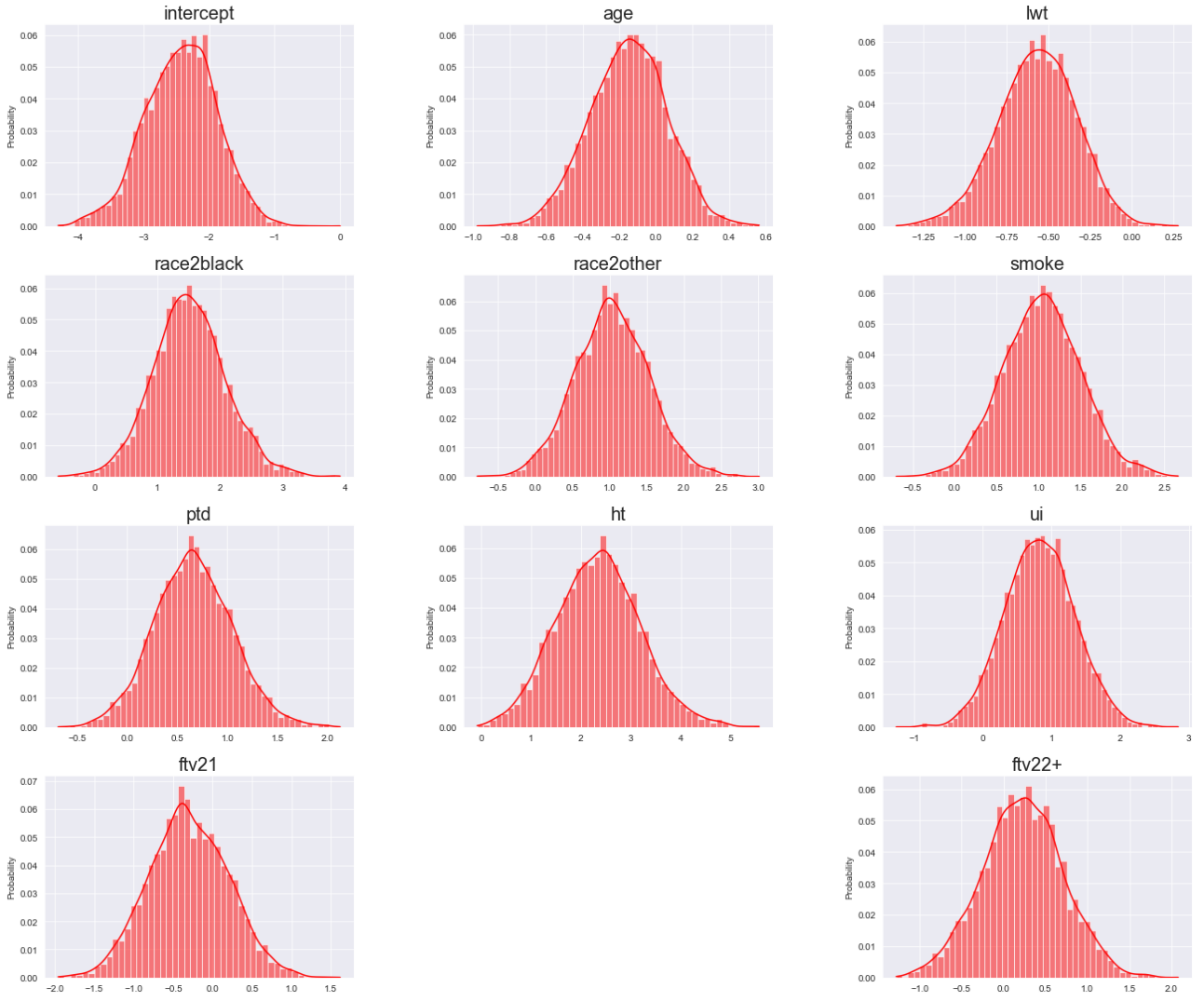


Figure 17: Histograms of the obtained coefficients

## References

- [1] Nawaf Bou-Rabee and Jesus Maria Sanz-Sern. *"Geometric integrators and the Hamiltonian Monte Carlo method"*. Acta Numerica 27, 2018.
- [2] Simon Duane, A.B. Kennedy, Bryan Pendleton, and Duncan Roweth. "hybrid monte carlo". pages 216–222, 1987.
- [3] Stanley Lemeshow Hosmer, David W and Rodney X Sturdivant. *"Applied Logistic Regression"*. John Wiley Sons, Inc., 1989.
- [4] Radford M. Neal. *"MCMC using Hamiltonian dynamics "*. Chapman Hall/ CRC Press, 2011.
- [5] Fabio Nobile. *"Lecture notes, Stochastic Simulation"*. 2022.
- [6] Samuel Thomas and Wanzhu Tu. "learning hamiltonian monte carlo in r". 2020.

## 9 Code

The following script presents the main functions used to get results in section 6.

```
1 import numpy as np
2 import scipy.stats as st
3 import matplotlib.pyplot as plt
4
5
6 def f(q, alpha): # unnormalized distribution
7     return np.exp(-alpha*(q@q-1/4)**2)
8
9 def U(q, alpha): #-log(f) (potential energy)
10    return alpha*(q@q-1/4)**2
11
12 def K(p, m): #kinetic energy
13    return np.sum(1/2*np.divide(p*p, m))
14
15 def dU(q, alpha): #gradient of U
16    dU1 = 4*q[0]*alpha*(q@q-1/4)
17    dU2 = 4*q[1]*alpha*(q@q-1/4)
18    return np.asarray([dU1, dU2])
19
20
21 def Verlet(q0, p0, eps, T, alpha, m): #Verlet's scheme
22    t=0
23    q=q0
24    p=p0
25    while t<T:
26        p_tmp = p - eps/2*dU(q, alpha)
27        q = q + eps*np.divide(p_tmp, m)
28        p = p_tmp - eps/2*dU(q, alpha)
29        t=t+eps
30    return q, p
31
32 def Hamiltonian_Monte_Carlo(q0, m, N, T, eps, alpha):
33    q = np.zeros([N+1, 2])
34    q[0, :] = q0
35    accepted = 0
36    rejected = 0
37    for i in range(N):
38        p = st.norm.rvs(loc=0, scale=np.sqrt(m))
39        q_star, p_star = Verlet(q[i, :], p, eps, T, alpha, m)
40        u = np.random.uniform()
41        if (u<np.exp(-U(q_star, alpha)+U(q[i, :], alpha)-K(p_star, m)+K(p, m)
42        ))):
43            q[i+1, :] = q_star
44            accepted = accepted+1
45        else:
46            q[i+1, :] = q[i, :]
47            rejected = rejected+1
48    ratio = accepted/(accepted+rejected)
49    return q, ratio
50
51 def Metropolis_Hastings(q0, N, alpha, sigma):
52    q = np.zeros([N+1, 2])
53    q[0, :] = q0
54    accepted = 0
55    rejected = 0
56    for i in range(N):
57        q_star = np.random.multivariate_normal(mean=q[i, :], cov=sigma*np.
58        eye(2))
```

```

57     a = f(q_star, alpha)/np.maximum(f(q[i, :], alpha), np.finfo(np.
float64).eps)
58     u = np.random.uniform()
59     if (u<a):
60         q[i+1, :] = q_star
61         accepted = accepted+1
62     else:
63         q[i+1, :] = q[i, :]
64         rejected = rejected+1
65     ratio = accepted/(accepted+rejected)
66     return q, ratio
67
68
69 def autocovariance(X): # Autocovariance of a Markov chain X, X tail of q
after burn-in lag B
70     mu = np.mean(X, axis=0)
71     N = len(X[:, 0])
72     cov = np.zeros([N-1, 2])
73
74     for k in range(N-1):
75         for i in range(N-k-1):
76             cov[k, :] += ((X[i+k, :]) - mu) * (X[i, :] - mu)
77             cov[k, :] = (1/(N-k-1)) * cov[k, :]
78
79     return cov
80
81 def autocorrelation(autocov): # Autocorrelation of a Markov chain,
calculated from its autocovariance
82     corr = np.zeros(np.shape(autocov))
83     for i in range(int(len(autocov[0, :, 0]))):
84         corr[:, i, 0] = np.divide(autocov[:, i, 0], np.maximum(np.abs(
autocov[:, 0, 0]), np.finfo(np.float64).eps))
85         corr[:, i, 1] = np.divide(autocov[:, i, 1], np.maximum(np.abs(
autocov[:, 0, 1]), np.finfo(np.float64).eps))
86
87     return corr
88
89 def get_M(autocov): # Last element of the chain used to calculate sigma^2
_mcmc
90     M1 = None
91     for k in range(int(len(autocov[:, 0])/2)):
92         if ((autocov[2*k, 0] + autocov[2*k+1, 0]) <= 0):
93             M1 = 2*k
94             break
95     if M1 is None:
96         M1 = int(len(autocov[:, 0]) - 2)
97     M2 = None
98     for j in range(int(len(autocov[:, 1])/2)):
99         if ((autocov[2*j, 1] + autocov[2*j+1, 1]) <= 0) :
100             M2 = 2*j
101             break
102     if M2 is None:
103         M2 = int(len(autocov[:, 1]) - 2)
104
105     return M1, M2
106
107 def get_sigma(autocov): # Time-average variance constant of a Markov chain,
calculated from its autocovariance
108     [M1, M2] = get_M(autocov)
109     id1 = range(1, M1)
110     S1 = autocov[0, 0] + 2*np.sum(autocov[id1, 0])
111     id2 = range(1, M2)

```

```

112     S2 = autocov[0, 1] + 2*np.sum(autocov[id2, 1])
113
114     return [S1, S2]
115
116 def effective_sample_size(autocov): # effective sample size of a Markov
chain, calculated from its autocovariance
117     [n, N] = np.shape(autocov[:, :, 0])
118     ess = np.zeros([n, 2])
119     for i in range(n):
120         sigma = get_sigma(autocov[i, :, :])
121         ess[i, 0] = N*np.divide(autocov[i, 0, 0], np.maximum(sigma[0], np.
finfo(np.float64).eps))
122         ess[i, 1] = N*np.divide(autocov[i, 0, 1], np.maximum(sigma[1], np.
finfo(np.float64).eps))
123
124     return ess

```

The following script presents the main functions used to get results in section 7

```

1 import seaborn as sns
2 import pandas as pd
3 import os
4 import scipy.stats as st
5 import matplotlib.pyplot as plt
6 import numpy as np
7
8
9 def log_posterior(q, X, y, sigma=1000): # log posterior
10
11     n, p = (X.shape[0], X.shape[1] - 1)
12
13     return q @ X.T @ (y.T - np.ones(n)) - np.ones(n) @ (
14         np.log1p(np.exp(-X @ q)) - q @ q.T / (2 * sigma)
15     )
16
17
18 def grad_log(X, y, q, sigma): # gradient of the log posterior
19     value = np.exp(-X @ q)
20     place = np.where(value == np.inf)[0]
21
22     value = value / (1 + value)
23     value[place] = 1.0
24
25     grad = X.T @ (y - np.ones(y.shape[0]) + value) - q / sigma
26
27     return grad
28
29
30 def K(p, m): # "kinetic energy"
31     return np.sum(1 / 2 * np.divide(p * p, m))
32
33
34 def Verlet(q0, m, p0, eps, T, X, y, sigma): # Verlet method
35     t = 0
36     q = q0
37     p = p0
38     while t < T:
39         p_tmp = p + eps / 2 * grad_log(X, y, q, sigma)
40
41         q = q + eps * np.divide(p_tmp, m)
42
43         p = p_tmp + eps / 2 * grad_log(X, y, q, sigma)
44

```

```

45     t = t + eps
46     return q, p
47
48
49 def Hamiltonian_Monte_Carlo(
50     q0, m, N, T, eps, X, y, sigma
51 ): # HMC scheme adapted for the low-weight problem
52
53     q = np.zeros([N + 1, q0.shape[1]])
54     q[0, :] = q0
55     accepted = 0
56     rejected = 0
57     for i in range(N):
58         p = st.norm.rvs(loc=0, scale=np.sqrt(m))
59
60         q_star, p_star = Verlet(
61             q0=q[i, :], m=m, p0=p, eps=eps, T=T, X=X, y=y, sigma=sigma
62         )
63
64         u = np.random.uniform()
65         if u < np.exp(
66             log_posterior(q_star, X, y, sigma)
67             - log_posterior(q[i, :], X, y, sigma)
68             - K(p_star, m)
69             + K(p, m)
70         ):
71             q[i + 1, :] = q_star
72             accepted = accepted + 1
73         else:
74             q[i + 1, :] = q[i, :]
75             rejected = rejected + 1
76             print(f"Problem:{i+1}")
77     ratio = accepted / (accepted + rejected)
78     return q, ratio
79
80
81
82 def MH_one_at_a_time(q0, N, var, X, y, sigma):
83     n = len(q0[0, :])
84     print(n)
85     q = np.zeros([N + 1, n])
86
87     q[0, :] = q0
88     accepted = 0
89     rejected = 0
90
91     for i in range(N):
92         j = np.random.randint(0, n)
93         q_star = np.copy(q[i, :])
94         # define posterior wrt other variates (or.....use the prior, like MH
95         ?)
96         q_star[j] = np.random.normal(loc=q[i, j], scale=var)
97         a = np.exp(log_posterior(q_star, X, y, sigma)-log_posterior(q[i, :],
98             X, y, sigma))
99         u = np.random.uniform()
100         if (u<a):
101             q[i+1, :] = q_star
102             accepted = accepted+1
103         else:
104             q[i+1, :] = q[i, :]
105             rejected = rejected+1

```



```

105     ratio = accepted/(accepted+rejected)
106     return q, ratio
107
108
109 def dataset_import():
110     this_dir = os.path.dirname(os.getcwd())
111
112     data_dir = os.path.join(this_dir, "src\\birthwt.csv")
113     df = pd.read_csv(data_dir)
114     y = df["low"].copy()
115     df = df.drop("Unnamed: 0", axis=1)
116     df = df.drop("low", axis=1)
117     df = df.drop("bwt", axis=1)
118     df.rename(columns={"race": "af_am"}, inplace=True)
119     df.insert(3, "Other", df["af_am"], True)
120     df.rename(columns={"ftv": "visit"}, inplace=True)
121     df.insert(9, "visits", df["visit"], True)
122
123     # transform to array
124     X = df.to_numpy()
125     y = y.to_numpy()
126     # visits last months
127     X[np.where(X[:, -1] < 2)[0], -1] = 0
128     X[np.where(X[:, -1] >= 2)[0], -1] = 1
129     # only one visit
130     X[np.where(X[:, -2] != 1)[0], -2] = 0
131     # afro american
132     X[np.where(X[:, 2] != 2)[0], 2] = 0
133     X[np.where(X[:, 2] == 2)[0], 2] = 1
134     # other
135     X[np.where(X[:, 3] != 3)[0], 3] = 0
136     X[np.where(X[:, 3] == 3)[0], 3] = 1
137
138     # intersept
139     X = np.c_[np.ones((X.shape[0], 1)), X]
140     # normalization of age and lwt, only non categorical variables
141     X[:, 1] = (X[:, 1] - np.mean(X[:, 1])) / np.std(X[:, 1])
142     X[:, 2] = (X[:, 2] - np.mean(X[:, 2])) / np.std(X[:, 2])
143     return X, y
144
145
146 def autocovariance(X): # X is only the tail of q after the burn-in
147     mu = np.mean(X, axis=0)
148     N = len(X[:, 0])
149     D = len(X[0, :])
150     cov = np.zeros([N-1, D])
151
152     for k in range(N-1):
153         for i in range(N-k-1):
154             cov[k, :] += ((X[i+k, :]) - mu) * (X[i, :] - mu)
155             cov[k, :] = (1/(N-k-1)) * cov[k, :]
156
157     return cov
158
159 def autocorrelation(autocov):
160     corr = np.zeros(np.shape(autocov))
161     D = len(autocov[0, 0, :])
162     for i in range(int(len(autocov[0, :, 0]))):
163         for k in range(D):
164             corr[:, i, k] = np.divide(autocov[:, i, k], np.maximum(np.abs(
autocov[:, 0, k]), np.finfo(np.float64).eps))
165

```

```

166     return corr
167
168 def get_M(autocov): # check adaptation of function for 2 dimensions
169     D = len(autocov[0, :])
170     M = np.zeros(D)
171     for d in range(D):
172         M[d] = None
173         for k in range(int(len(autocov[:, d])/2)):
174             if ((autocov[2*k, d] + autocov[2*k+1, d])<=0):
175                 M[d] = 2*k
176                 break
177         if M[d] is None:
178             M[d] = int(len(autocov[:, d])-2)
179
180
181     return M
182
183 def get_sigma(autocov): # To check (difference between serie 13 and lecture
184     # notes)
185     M = np.ndarray.astype(get_M(autocov), int)
186     D = len(autocov[0, :])
187     S = np.zeros(D)
188     for d in range(D):
189         id = range(1, M[d])
190         S[d] = autocov[0, d] + 2*np.sum(autocov[id, d])
191
192     return S
193
194 def effective_sample_size(autocov):
195     [n, N, D] = np.shape(autocov[:, :, :])
196     ess = np.zeros([n, D])
197     for i in range(n):
198         sigma = get_sigma(autocov[i, :, :])
199         for d in range(D):
200             #print("Sigma=", sigma, ", c0=", autocov[i, 0, :], ", N=", N)
201             ess[i, d] = N*np.divide(autocov[i, 0, d], np.maximum(sigma[d],
202             np.finfo(np.float64).eps))
203
204     return ess
205
206 if __name__ == "__main__":
207
208     X, y = dataset_import()
209
210     q0 = np.zeros((1, X.shape[1]))
211     eps = 0.01
212     m = np.ones(X.shape[1])
213     T = 0.40
214     N = 6000
215     sigma = 1000
216     B = 4800
217
218     (q, ratio) = Hamiltonian_Monte_Carlo(q0, m, N, T, eps, X, y, sigma)
219
220     # histograms
221     sns.set_style("darkgrid")
222     rel = sns.histplot(
223         data=q[:, 0],
224         kde=True,
225         stat="probability",

```

```

226         color="r",
227         label="probabilities",
228     )
229     plt.title("intercept", fontsize=20)
230     plt.show()
231
232     rel = sns.histplot(
233         data=q[:, 1],
234         kde=True,
235         stat="probability",
236         color="r",
237         label="probabilities",
238     )
239     plt.title("age", fontsize=20)
240     plt.show()
241
242     rel = sns.histplot(
243         data=q[:, 2],
244         kde=True,
245         stat="probability",
246         color="r",
247         label="probabilities",
248     )
249     plt.title("lwt", fontsize=20)
250     plt.show()
251
252     rel = sns.histplot(
253         data=q[:, 3],
254         kde=True,
255         stat="probability",
256         color="r",
257         label="probabilities",
258     )
259     plt.title("race2black", fontsize=20)
260     plt.show()
261
262     rel = sns.histplot(
263         data=q[:, 4],
264         kde=True,
265         stat="probability",
266         color="r",
267         label="probabilities",
268     )
269     plt.title("race2other", fontsize=20)
270     plt.show()
271
272     rel = sns.histplot(
273         data=q[:, 5],
274         kde=True,
275         stat="probability",
276         color="r",
277         label="probabilities",
278     )
279     plt.title("smoke", fontsize=20)
280     plt.show()
281
282     rel = sns.histplot(
283         data=q[:, 6],
284         kde=True,
285         stat="probability",
286         color="r",
287         label="probabilities",

```

```

288 )
289 plt.title("ptd", fontsize=20)
290 plt.show()
291
292 rel = sns.histplot(
293     data=q[:, 7],
294     kde=True,
295     stat="probability",
296     color="r",
297     label="probabilities",
298 )
299 plt.title("ht", fontsize=20)
300 plt.show()
301
302 rel = sns.histplot(
303     data=q[:, 8],
304     kde=True,
305     stat="probability",
306     color="r",
307     label="probabilities",
308 )
309 plt.title("ui", fontsize=20)
310 plt.show()
311
312 rel = sns.histplot(
313     data=q[:, 9],
314     kde=True,
315     stat="probability",
316     color="r",
317     label="probabilities",
318 )
319 plt.title("ftv21", fontsize=20)
320 plt.show()
321
322 rel = sns.histplot(
323     data=q[:, 10],
324     kde=True,
325     stat="probability",
326     color="r",
327     label="probabilities",
328 )
329 plt.title("ftv22+", fontsize=20)
330 plt.show()

```