

Lecture 3: Markov Chain Monte Carlo

Ziyu Shao

School of Information Science and Technology
ShanghaiTech University

March 19, 2025

Outline

- 1 Introduction of MCMC
- 2 Metropolis–Hastings Algorithm
- 3 Gibbs Sampler
- 4 Hamiltonian Monte Carlo Method
- 5 References

Outline

- 1 Introduction of MCMC
- 2 Metropolis–Hastings Algorithm
- 3 Gibbs Sampler
- 4 Hamiltonian Monte Carlo Method
- 5 References

Markov Chain Monte Carlo (MCMC)

- Revolutionized statistics and scientific computation
- Expanding the range of possible distributions that we can simulate from, including joint distributions in high dimensions
- Basic idea: build your own Markov chain (X_0, X_1, \dots) so that the desired distribution π is the stationary distribution of the chain.
- Sampling from distribution π (running the chain for a long time and then sampling)
- Further do sample mean & sample variance & other sample functions
- Connections to optimization

Sampling ← optimization

MCMC Method

- Forward engineering(Analysis): given the transition matrix P , find the stationary distribution of Markov chain.
- Reverse engineering(Design): given a distribution π that we want to simulate, we will engineer a Markov chain whose stationary distribution is π . Then run this engineered Markov chain for a long time, the distribution of the chain will approach π .

MCMC Method

Non-reversible M.C.

- Markov Chain Monte Carlo (MCMC) is a remarkable methodology, which utilizes Markov sequences to effectively simulate from what would otherwise be intractable distributions.
- All MCMC algorithms construct reversible (time-reversible) Markov chain: detailed balance equations help us.
- Two most widely used algorithms: Metropolis-Hastings & Gibbs Sampling

Theory Justification: Strong Law of Large Numbers for Markov Chains

Theorem

Assume that X_0, X_1, \dots is an ergodic Markov chain with stationary distribution π . Let r be a bounded, real-valued function. Let X be a random variable with distribution π . Then, with probability 1,

$$\lim_{n \rightarrow \infty} \frac{r(X_1) + \dots + r(X_n)}{n} = E(r(\underline{\underline{X}})).$$

where $E(r(X)) = \sum_j r(j)\pi_j$.

Example: Bob's Lunch

Stationary π :

distribution

Bob's daily lunch choices at the cafeteria are described by a Markov chain with transition matrix

$$P = \begin{pmatrix} \text{Yogurt} & 1 & \text{Yogurt} & 1 \\ \text{Salad} & 2 & \text{Salad} & 2 \\ \text{Hamburger} & 3 & \text{Hamburger} & 3 \\ \text{Pizza} & 4 & \text{Pizza} & 4 \end{pmatrix}$$

where the columns represent the current state and the rows represent the next state.

Yogurt costs \$3.00, hamburgers cost \$7.00, and salad and pizza cost \$4.00 each. Over the long term, how much, on average, does Bob spend for lunch?

$$\pi P = \pi$$

$$\pi_1 = \pi_2 = \frac{1}{5},$$

$$\pi_3 = \frac{3}{20};$$

$$\pi_4 = \frac{9}{20};$$

Solution

1^o. $X = \text{state of M.C.} ; X \in \{ \text{"Yogurt"}, \text{"Salad"}, \text{"Hamburger"} \}$

$$r(x) : \text{food rice} \quad r(x) = \begin{cases} 3 & \text{if } x = \text{"Yogurt"} \\ 4 & \text{if } x = \text{"Salad" or "Pizza"} \\ 7 & \text{if } x = \text{"Hamburger"} \end{cases}$$

2^o. n : days for lunch or M.C. $\{X_t\}$. Ergodic (irreducible;
aperiodic;
finite state;)

$$\lim_{n \rightarrow \infty} \frac{1}{n} [r(X_1) + r(X_2) + \dots + r(X_n)] \stackrel{\text{SLN}}{=} E[r(X)] = \sum_x r(x) \pi(x)$$

3^o. Stationary distribution of M.C. $\pi p = \pi$; $\pi = (\frac{1}{5}, \frac{1}{5}, \frac{3}{20}, \frac{9}{20})$

$$\sum_x n x \pi(x) = 3 \cdot \frac{1}{5} + 4 \cdot \frac{1}{5} + 7 \cdot \frac{3}{20} + 4 \cdot \frac{9}{20} \\ = 4.25 \text{ $ per day}$$

Example: Binary Sequences with No Adjacent 1s

$m=2$; therefore $2^M = 2^2 = 4$ binary sequences.

$\begin{array}{c} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{array}$	$\begin{array}{c} 00 \\ 01 \\ 10 \\ 11 \end{array}$	\checkmark	good sequences.
$\begin{array}{c} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{array}$	$\begin{array}{c} X \\ X \end{array}$	\times	

$$\frac{1}{3} \cdot 0 + \frac{1}{3} \cdot 1 + \frac{1}{3} \cdot 1 = \frac{2}{3}$$

$m=100$; $2^{100} \approx 10^{30}$ binary sequences; $\approx 10^{21}$ good sequences.

Consider sequences of length m consisting of 0s and 1s. Call a sequence good if it has no adjacent 1s. What is the expected number of 1s in a good sequence if all good sequences are equal likely?

1°. Construct an ergodic M.c. $\left\{ \begin{array}{l} \text{State : good sequence} \\ \text{state space : good sequences} \end{array} \right.$

desired stationary distribution π : uniform distribution.

2°: X : state (good sequences).

$r(x)$: # of "1" in state X .

[# of good sequences
= 0]

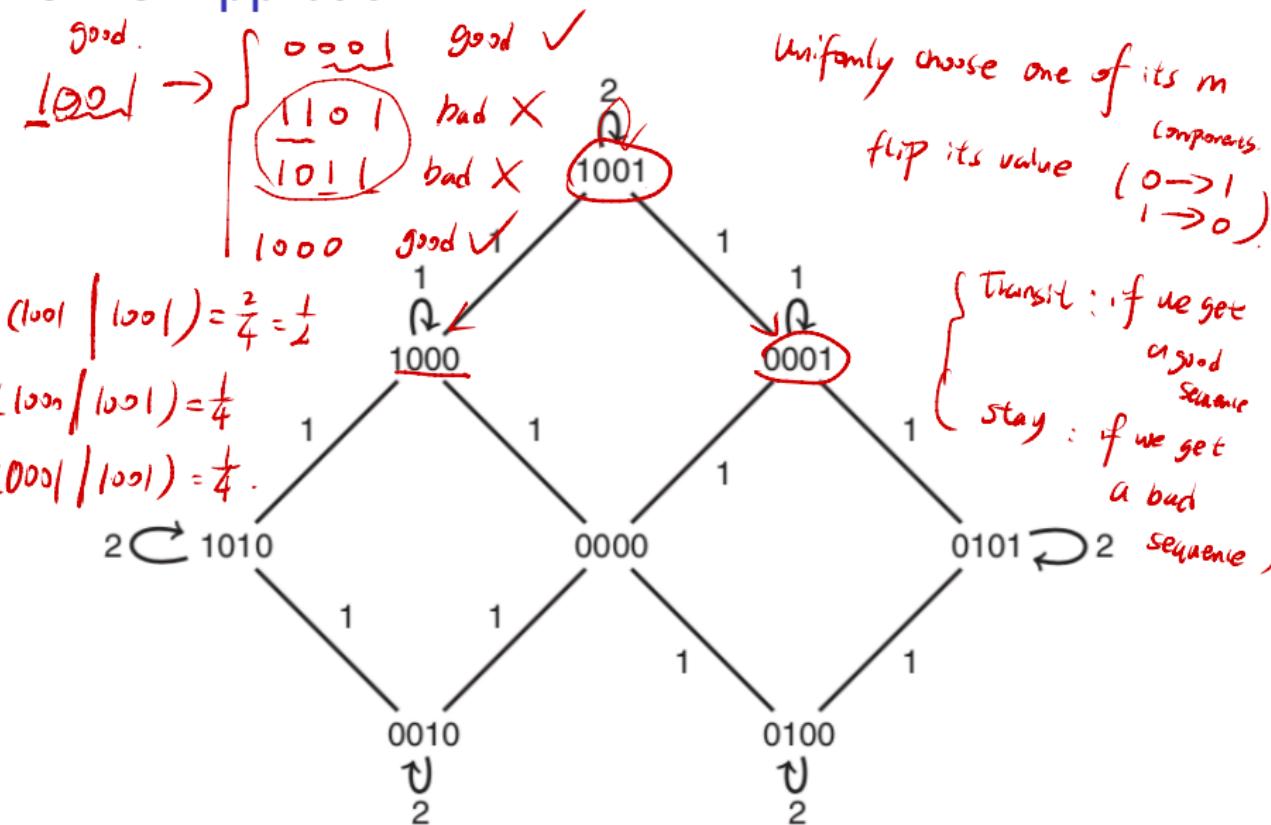
$\left[\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \dots, \frac{1}{2} \right]$

$$E[r(X)] \approx \frac{1}{n} [r(x_1) + r(x_2) + \dots + r(x_n)] \text{ finite}_n$$

MCMC Approach

$m = 4$

Design P (transition matrix)
prob.



MCMC Approach

1^o. DTM C

finite state

$$< 2^m$$

ergodic

irreducible

a priori: c.

$$\textcircled{a}$$

length m

$$\textcircled{a}$$

$$0000 \rightarrow \textcircled{b}$$

$$\textcircled{b} \rightarrow \textcircled{b}$$

at most

$$m-1 \\ \text{flip}$$

at most

$$m-1 \\ \text{flip}$$

$$a \rightleftharpoons b$$

$\leq 2(m-1)$ steps

2^o. $\textcircled{a} \rightarrow \textcircled{b}$. iff $\text{dist}^H(a, b) = 1$

$$P_{a,b} = \frac{1}{m}$$

$$\textcircled{a \neq b}$$

Hamming
distance

$$\begin{array}{r} 0000 \\ 1000 \\ \hline \end{array} \quad \text{dist}^H(\cdot, \cdot) = 1$$

$$\begin{array}{r} 0001 \\ 1000 \\ \hline \end{array} \quad \text{dist}^H(\cdot, \cdot) = 2$$

Solve detailed balance equation DBE.

$$\pi_i P_{i,j} = \pi_j P_{j,i}, \quad i \neq j, \quad \text{dist}^H(i, j) = 1$$

C: # of
good moves

$$P_{i,j} = P_{j,i} = \frac{1}{m} \Rightarrow \pi_i = \pi_j \quad \forall i \neq j \Rightarrow \pi_i = \frac{1}{C}$$

Solution of DBE: uniform distribution (π) \Rightarrow stationary distribution of M.C.

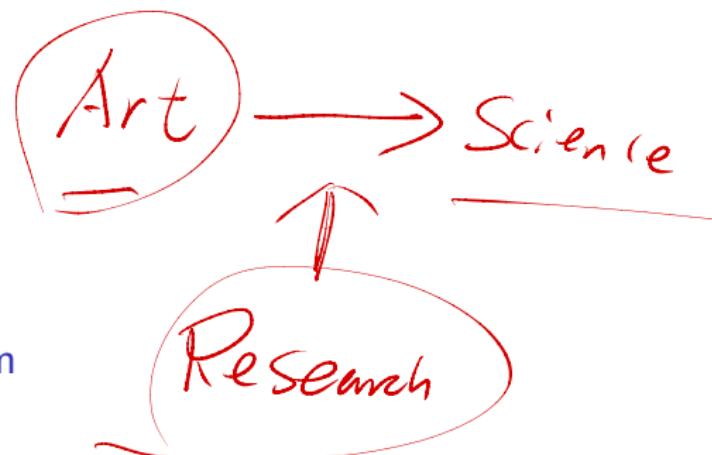
3^o. $m = 100$; theory $H = 27.7921$. Sample $n = 100000$; MCMC $\bar{H} = 27.833$.

HW: CTMC?

MCMC Approach

Outline

- 1 Introduction of MCMC
- 2 Metropolis–Hastings Algorithm
- 3 Gibbs Sampler
- 4 Hamiltonian Monte Carlo Method
- 5 References



Basic Idea

- Proposed by Nicholas Metropolis in 1953 & further developed by Wilfred Keith Hastings in 1970.
- Start with any irreducible Markov chain on the state space of interest
- Then modify it into a new Markov chain with desired stationary distribution
- Modification: moves are proposed according to the original chain, but the proposal may or may not be accepted.
- Art. choice of the probability of accepting the proposal

proposal M.C.

Algorithm for DTMC

Let $\pi = (\pi_1, \dots, \pi_M)$ be a desired stationary distribution on state space $\{1, \dots, M\}$. Assume that $\pi_i > 0$ for all i (if not, just delete any states i with $s_i = 0$ from the state space). Suppose that $P = (p_{ij})$ is the transition matrix for a Markov chain on state space $\{1, \dots, M\}$. Intuitively, P is a Markov chain that we know how to run but that doesn't have the desired stationary distribution. Now we will modify P to construct a Markov chain.

- Use the original transition probabilities p_{ij} to propose where to go next
- Then accept the proposal with probability a_{ij}
- Staying in the current state in the event of a rejection.
- Normalization of π does not need to be known.

Algorithm Metropolis-Hastings

Require:

Stationary distribution $\pi = (\pi_1, \dots, \pi_M)$;

Original transition matrix $P = (p_{ij})$; proposal M.C.

State X_0 (chosen randomly or deterministically);

Ensure:

Modified transition matrix $P' = (p'_{ij})$;

- 1: **repeat**
- 2: If $X_n = i$, propose a new state j using the transition probabilities in the i th row of the original transition matrix P ;
- 3: Compute the acceptance probability $a_{ij} = \min\left(\frac{\pi_j p_{ji}}{\pi_i p_{ij}}, 1\right)$;
- 4: Flip a coin that lands Heads with probability a_{ij} ;
- 5: If the coin lands Heads, accept the proposal (i.e., go to j), setting $X_{n+1} = j$. Otherwise, reject the proposal (i.e., stay at i), setting $X_{n+1} = i$;
- 6: **until** Convergence;
- 7: **return** P' ;

Metropolis–Hastings Algorithm

1^o. irreducible ✓ (from irreducible proposal M.C.)
aperiodic ✓

2^o. detailed balance equation,

$$\pi_i p_{i,j} \stackrel{?}{=} \pi_j p_{j,i} \quad | \quad p_{i,j}' = p_{i,j} \cdot a_{i,j}$$

Theorem

The sequence X_0, X_1, \dots constructed by the Metropolis–Hastings algorithm is a reversible Markov chain with stationary distribution π

$$\begin{aligned} \text{LHS: } \pi_i p_{i,j}' &= \pi_i \cdot p_{i,j} \underbrace{a_{i,j}}_{\min(\pi_j p_{j,i}, 1)} = \underbrace{\pi_i p_{i,j}}_{\min(\pi_j p_{j,i}, \pi_i p_{i,j})} \cdot \min\left(\frac{\pi_j p_{j,i}}{\pi_i p_{i,j}}, 1\right) \\ &= \min(\pi_j p_{j,i}, \pi_i p_{i,j}) \quad \min(A, B) \end{aligned}$$

$$\begin{aligned} \text{RHS: } \pi_j p_{j,i}' &= \pi_j p_{j,i} \cdot a_{j,i} = \underbrace{\pi_j p_{j,i}}_{\min(\pi_i p_{i,j}, \pi_j p_{j,i})} \cdot \min\left(\frac{\pi_i p_{i,j}}{\pi_j p_{j,i}}, 1\right) \\ &= \min(\pi_i p_{i,j}, \pi_j p_{j,i}) \quad \min(B, A) \end{aligned}$$

Remarks

$$\pi_x \propto e^{\beta \phi(x)}$$

$$\frac{\pi_x}{\pi_y} = e^{B[\phi(x) - \phi(y)]}$$

$$\pi_x = \frac{e^{\beta \phi(x)}}{\sum_y e^{\beta \phi(y)}} \rightarrow C$$

- The exact form of π is not necessary to implement Metropolis–Hastings. The algorithm only uses ratios of the form $\frac{\pi_j}{\pi_i}$. Thus, π needs only to be specified up to proportionality.
- If the proposal transition matrix P is symmetric,
 $a_{ij} = \min\left(\frac{\pi_j}{\pi_i}, 1\right)$. $p_{i,j} = p_{j,i}$
- The algorithm works for any irreducible proposal chain. Thus, the user has wide latitude to find a proposal chain that is efficient in the context of their problem.
- Under some conditions (e.g. finite state), the resulting Metropolis–Hastings chain is also ergodic with limiting(stationary) distribution π .

Remarks

- The generated sequence X_0, X_1, \dots, X_n gives an approximate sample from π .
- If the chain requires many steps to get close to stationarity, there may be initial bias.
- Burn-in: the practice of discarding the initial iterations and retaining X_m, X_{m+1}, \dots, X_n , for some m .
- The strong laws of large numbers for Markov chains:

$$\lim_{n \rightarrow \infty} \frac{r(X_m) + \dots + r(X_n)}{n - m + 1} = E(r(X)) = \sum_x r(x)\pi_x.$$

Example: Power-law Distribution

CTMC

Power-law distributions are positive probability distributions of the form $\pi_i \propto i^S$, for some constant S . Unlike distributions with exponentially decaying tails (e.g., Poisson, geometric, exponential, normal), power-law distributions have fat tails, and thus are often used to model skewed data. Let

$$\pi_i = \frac{i^{-3/2}}{\sum_{k=1}^{\infty} k^{-3/2}}, \text{ for } i = 1, 2, \dots$$

Implement a Metropolis–Hastings algorithm to simulate from π .

Solution 1°. State Space : $\{1, 2, \dots, \infty\}$

State = positive integer

Proposal Markov chain :

Simple Symmetric random walk

on the positive integer.

with one reflecting bound.



Irreducible M.C.

$$p_{i,j} = \begin{cases} \frac{1}{2} & \text{if } j = i \pm 1, i \geq 2 \\ 1 & \text{if } i = 1, j = 2 \\ 0 & \text{otherwise.} \end{cases}$$

2°. $\pi_i \propto i^{-\frac{3}{2}}$

the acceptance prob.

$$\begin{aligned} a_{i,j} &= \min\left(\frac{\pi_j p_{j,i}}{\pi_i p_{i,j}}, 1\right) \\ &= \min\left(\frac{j^{-\frac{3}{2}} \cdot p_{j,i}}{i^{-\frac{3}{2}} \cdot p_{i,j}}, 1\right) \end{aligned}$$

Solution $a_{1,2} = \min\left(\frac{x_2 p_{2,1}}{\underline{x_1 p_{1,2}}}, 1\right) = \min\left(\frac{2^{-\frac{3}{2}} \cdot 2^{-1}}{1^{-\frac{1}{2}} \cdot 1}, 1\right)$
 $= 2^{-\frac{5}{2}}$

$a_{2,1} = \min(2^{\frac{5}{2}}, 1) = 1$

3°. $i, j \geq 2$ $\Rightarrow a_{i,i+1} = \min\left(\frac{(i+1)^{-\frac{3}{2}}}{i^{-\frac{3}{2}}}, 1\right)$
 $\underline{p_{i,j} = p_{j,i} = \frac{1}{2}}$
 $= \left(\frac{i}{i+1}\right)^{\frac{3}{2}}$

$a_{i+1,i} = \min\left(\frac{i^{-\frac{3}{2}}}{(i+1)^{-\frac{3}{2}}}, 1\right)$
 $= \min\left(\left(\frac{i+1}{i}\right)^{\frac{3}{2}}, 1\right) = 1$

Solution

TABLE 5.1 Comparison of Markov chain Monte Carlo Estimates with Exact Probabilities for Power-Law Distribution

i	1	2	3	4	5	6	7	8	≥ 9
Simulation	0.389	0.137	0.075	0.048	0.034	0.026	0.021	0.017	0.252
Exact	0.383	0.135	0.074	0.048	0.034	0.026	0.021	0.017	0.262

Example: Zipf Distribution Simulation

Let $M \geq 2$ be an integer. An r.v. X has the *Zipf distribution* with parameter $a > 0$ if its PMF is

$$P(X = k) = \frac{1/k^a}{\sum_{j=1}^M (1/j^a)},$$

for $k = 1, 2, \dots, M$ (and 0 otherwise). This distribution is widely used in linguistics for studying frequencies of words.

Create a Markov chain X_0, X_1, \dots whose stationary distribution is the Zipf distribution, and such that $|X_{n+1} - X_n| \leq 1$ for all n . Your answer should provide a simple, precise description of how each move of the chain is obtained, i.e., how to transition from X_n to X_{n+1} for each n .

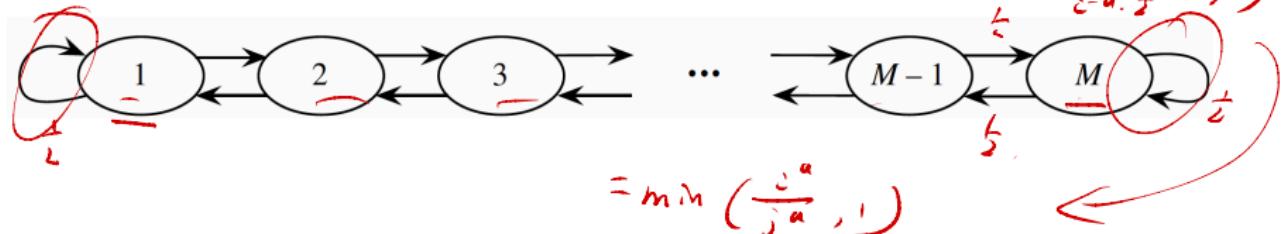
Solution

1°. Proposal Markov chain : state $i, 1 \leq i \leq n$

Statespace : $\{1, 2, \dots, n\}$

Simple random walk : $p_{i,j} = \begin{cases} \frac{1}{2} & \text{iff } |i-j| \leq 1 \\ 0 & \text{otherwise.} \end{cases}$

$$2^{\circ}. |i-j| \leq 1, a_{i,j} = \min \left(\frac{x_j p_{i,j}}{x_i p_{j,i}}, 1 \right) = \min \left(\frac{(j-a) \cdot \frac{1}{2}}{(i-a) \cdot \frac{1}{2}}, 1 \right)$$



$$a_{1,1} = a_{M,M}$$

$$= 1$$

$$a_{i,i+1} = \min \left(\frac{c^a}{(i+1)^a}, 1 \right) = \left(\frac{c}{i+1} \right)^a, 1 \leq i \leq M-1$$

$$a_{i,i-1} = \min \left(\left(\frac{c}{i-1} \right)^a, 1 \right) = 1, 2 \leq i \leq M$$

Solution

Example: Knapsack Problem

Sampling → optimization

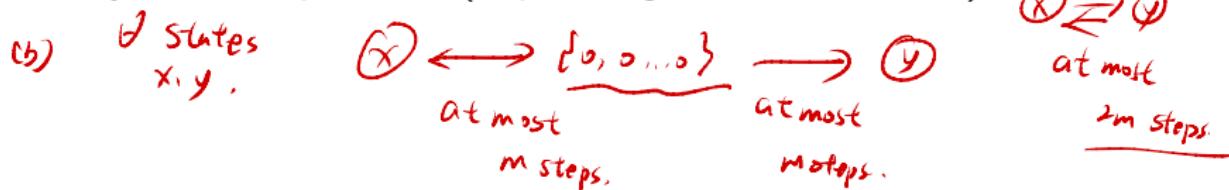
- m treasures with labels from 1 to m , where the j th treasure is worth g_j gold pieces and weighs w_j pounds.
- The maximum weight we can carry is w pounds.
- We must choose a vector $x = (x_1, \dots, x_m)$, where x_j is 1 if we choose the j th treasure and 0 otherwise, such that the total weight of the treasures j with $x_j = 1$ is at most w .
- Let C be the space of all such vectors, so C consists of all binary vectors (x_1, \dots, x_m) with $\sum_{j=1}^m x_j w_j \leq w$.
- We wish to maximize the total worth of the treasure we take.
- Finding an optimal solution is an extremely difficult problem, known as the knapsack problem, which has a long history in computer science.
- A brute force solution would be completely infeasible in general. How about MCMC?

Example: Knapsack Problem

(a) State $x = (x_1, \dots, x_m)$; State space $C = \{x : \sum_{j=1}^m x_j w_j \leq w\}$

if $x \rightarrow y$; $P_{x,y} = \frac{1}{m}$; $\pi_P = \pi$ ✓ π (uniform distribution)
 $y \rightarrow x$; $P_{y,x} = \frac{1}{m}$;

- (a) Consider the following Markov chain. Start at $(0, 0, \dots, 0)$. One move of the chain is as follows. Suppose the current state is $x = (x_1, \dots, x_m)$. Choose a uniformly random j in $\{1, 2, \dots, m\}$, and obtain y from x by replacing x_j with $1 - x_j$ (i.e., toggle whether that treasure will be taken). If y is not in C , stay at x ; if y is in C , move to y . Show that the uniform distribution over C is stationary for this chain.
- (b) Show that the chain from (a) is irreducible, and that it may or may not be aperiodic (depending on w, w_1, \dots, w_m).



Solution

periodicity : w, w_1, \dots, w_n

1^o . $w_1 + w_2 + \dots + w_m < w$

All binary vectors of length m are allowed.

$(0, 0, \dots, 0)$ period of \Rightarrow (pick up + drop)

2^o . $w_1 = w$

$(0, x_2, \dots, x_m)$ \otimes



$\frac{1}{m} \rightarrow \underline{(0, x_2, \dots, x_m)} \quad \textcircled{Y} \notin C.$

aperiodic

Solution

$$\max_{i=1,\dots,n} x_i \approx \log\text{-sum-exp.}$$

$$\frac{1}{\beta} \log \left(\sum_{i=1}^n e^{\beta x_i} \right)$$

✓ Conjugate-transform

Entropy regularization

Example: Knapsack Problem

$$\max_{x \in C} V(x) \iff \max_{x \in C} z(x) V(x)$$

x^* is unique.

- (c) The chain from (a) is a useful way to get approximately uniform solutions, but Bilbo is more interested in finding solutions where the value (in gold pieces) is high. In this part, the goal is to construct a Markov chain with a stationary distribution that puts much higher probability on any particular high-value solution than on any particular low-value solution. Specifically, suppose that we want to simulate from the distribution

$$\frac{\max_{x \in C} z(x) V(x)}{\sum_{x \in C} z(x)} = \frac{e^{\beta V(x)}}{\sum_y e^{\beta V(y)}}$$

$s(x) \propto e^{\beta V(x)}$

$\beta > 0$

where $V(x) = \sum_{j=1}^m x_j g_j$ is the value of x in gold pieces and β is a positive constant. The idea behind this distribution is to give exponentially more probability to each high-value solution than to each low-value solution. Create a Markov chain whose stationary distribution is as desired.

Solution

(B20)

1^o. desired distribution $\pi(x) \propto e^{\beta V(x)}$ - $V(x) = \sum_{j=1}^n x_j \cdot g_j$

2^o. Proposal Markov chain in part(a)
Acceptance prob.

$$\alpha_{x,y} = \min\left(\frac{z_y p_{y|x}}{z_x p_{x|y}}, 1\right)$$

$$= \min\left(\frac{z_y}{z_x}, 1\right) = \min\left[e^{\beta(V(y)-V(x))}, 1\right]$$

$$\pi(x) \propto (e^{\beta V(x)})$$

(1)

$$\pi(x) \propto (F[V(x)])$$

Solution

$$\pi(x) \propto e^{\beta V(x)} \Rightarrow \pi(x) = \frac{e^{\beta V(x)}}{\sum_{y \in C} e^{\beta V(y)}}, \forall x \in C.$$

1^o. $\beta \rightarrow 0$: $\pi(x) \rightarrow \frac{1}{|C|}$, $|C|$ is size of C .

Uniform distribution.

(c) We can apply Metropolis-Hastings using the chain from (a) to make proposals. Start at $(0, 0, \dots, 0)$. Suppose the current state is $x = (x_1, \dots, x_m)$. Then:

1. Choose a uniformly random J in $\{1, 2, \dots, m\}$, and obtain y from x by replacing x_J with $1 - x_J$.

2. If y is not in C , stay at x . If y is in C , flip a coin that lands Heads with probability $\min(1, e^{\beta(V(y) - V(x))})$. If the coin lands Heads, go to y ; otherwise, stay at x .

2^o. $\beta \rightarrow \infty$; $V(x^*) > V(x)$, $\forall x \neq x^*$ $\pi(x) \rightarrow \begin{cases} 1 & \text{if } x = x^* \\ 0 & \text{otherwise.} \end{cases}$

$$\pi(x^*) = \frac{e^{\beta V(x^*)}}{\sum_{y \neq x^*} e^{\beta V(y)} + e^{\beta V(x^*)}} = \frac{1}{\sum_{y \neq x^*} e^{\beta(V(y) - V(x^*))} + 1} \rightarrow \frac{1}{0+1} = 1$$

$$\pi(x) \rightarrow 0, x \neq x^*,$$



Solution $a_{x,y} = \min(1, e^{\beta [v(y) - v(x)]})$.

1^o. $\beta > 1$. if x is local optimal , y : neighbor states of x .

Bruce Hajek 1988.

$$T_t = \frac{c}{\log(t+1)}$$

$$\beta_t = \log(t+1)$$

$$\textcircled{B \rightarrow \infty} \quad a_{x,y} = 0,$$

$$v(y) < v(x)$$

stuck in local optimal.

slow convergence

} Exploration

2^o. $\beta \ll 1$. $\textcircled{B \rightarrow 0}$, $a_{x,y} = 1$

far from optimal.

Fast Convergence

} Exploration

β :
Small
Large
explore state space

$$\textcircled{B = \frac{1}{Temperature}}$$

T_t

Annealing
Simulated

Discrete Time, Continuous State Space

- MH algorithm can also be used in the discrete time, continuous state case (Discrete Markov Process), when π is a probability density function.
- For a continuous state space S , a transition *function* replaces the transition matrix.
- $f_{x,y} = f_{X_{n+1}|X_n}(y|x)$ is the one-step transition probability density from state x to state y .
- Irreducible condition: $f_{x,y} > 0, \forall x, y \in S$
- Compute the acceptance probability:

$$a_{i,j} = \min \left(\frac{\pi_j f_{j,i}}{\pi_i f_{i,j}}, 1 \right).$$

- π_j : pdf in state i
- $f_{i,j}$: one-step transition probability density from state i to state j

Example: Beta Simulation

1^o. PDF $f(x) = \frac{1}{\beta(a,b)} x^{a-1} (1-x)^{b-1}, 0 < x < 1$
 $a > 0, b > 0$

Suppose that we want to generate $W \sim \text{Beta}(a, b)$. What we have available are i.i.d. $\text{Unif}(0, 1)$ r.v.s. How can we generate W which is *approximately* $\text{Beta}(a, b)$ if a and b are any positive real numbers, with the help of a Markov chain on the continuous state space $(0, 1)$?

Hw: Acceptance - Rejection vs. MCMC(MH)
Method

2^o.

Solution

Solution

1^o. Proposal M.P.: independent sampler. T degeneration of M.p)

Unif(0,1) r.v.

2^o. Acceptance prob. $\min_{w \rightarrow u} \left(\frac{z_u \cdot f_w}{z_w \cdot f_u}, 1 \right)$

Let W_0 be any starting state, and generate a chain W_0, W_1, \dots as follows. If the chain is currently at state w (a real number in $(0, 1)$), then:

1. Generate a proposal u by drawing a Unif(0, 1) r.v.

2. Accept the proposal with probability $\min \left(\frac{u^{a-1}(1-u)^{b-1}}{w^{a-1}(1-w)^{b-1}}, 1 \right)$. If the proposal is accepted, go to u ; otherwise, stay at w .

$$= \min \left(\frac{z_u}{z_w}, 1 \right)$$

$$= \min \left(\frac{u^{a-1}(1-u)^{b-1}}{w^{a-1}(1-w)^{b-1}}, 1 \right)$$

Outline

- 1 Introduction of MCMC
- 2 Metropolis–Hastings Algorithm
- 3 Gibbs Sampler
- 4 Hamiltonian Monte Carlo Method
- 5 References

Basic Idea

- Proposed by brothers Stuart and Donald Geman in 1984.
- Named after the physicist Josiah Willard Gibbs, in reference to an analogy between the sampling algorithm and statistical physics.
- Obtaining approximate draws from a joint distribution, based on sampling from conditional distributions one at a time
- Especially useful in problems where these conditional distributions are pleasant to work with.

Basic Idea

$$X = (\underline{x}_1, \dots, \underline{x}_n)$$

Joint distribution

$t \leftarrow 1$

Systematic Scan

$$\left\{ \begin{array}{l} f(x_1^t | x_2^{t-1}, \dots, x_m^{t-1}) \\ f(x_2^t | \underline{x}_1^t, \underline{x}_3^{t-1}, \dots, \underline{x}_m^{t-1}) \\ f(x_3^t | x_1^t, x_2^t, x_4^{t+1}, \dots, x_m^{t+1}) \\ f(x_m^t | x_1^t, x_2^t, \dots, x_{m-1}^t) \end{array} \right.$$

- At each stage, one variable is updated (keeping all the other variables fixed) by drawing from the conditional distribution of that variable given all the other variables. $t \leftarrow t + 1$

- Two major kinds of Gibbs sampler:

- ▶ systematic scan: the updates sweep through the components in a deterministic order.
- ▶ random scan: a randomly chosen component is updated at each stage.

$$f(x_2^t | x_1^{t-1}, x_3^{t+1}, \dots, x_m^{t-1})$$


Algorithm: Systematic Scan Gibbs Sampler

Let X and Y be discrete r.v.s with joint PMF

$p_{X,Y}(x,y) = P(X=x, Y=y)$. We wish to construct a two-dimensional Markov chain (X_n, Y_n) whose stationary distribution is $p_{X,Y}$. The systematic scan Gibbs sampler proceeds by updating the X -component and the Y -component in alternation. If the current state is $(X_n, Y_n) = (x_n, y_n)$, then we update the X -component while holding the Y -component fixed, and then update the Y -component while holding the X -component fixed.

Algorithm Systematic Scan Gibbs Sampler

Require:

Joint PMF $p_{X,Y}$;

Initial state (X_0, Y_0) ;

Ensure:

Two-dimensional Markov chain (X_n, Y_n) ;

- 1: **repeat**
- 2: Draw a value x_{n+1} from the conditional distribution $\underline{P(X|Y=y_n)}$ of X given $Y = y_n$, and set $X_{n+1} = x_{n+1}$;
- 3: Draw a value y_{n+1} from the conditional distribution of Y given $X = x_{n+1}$, and set $Y_{n+1} = y_{n+1}$;
- 4: **return** (X_{n+1}, Y_{n+1}) ;
- 5: **until** $n \geq N$;

$$\begin{aligned} & \underline{P(X|Y=y_n)} \\ & \cancel{P(Y|X=x_n)} \end{aligned}$$

Algorithm: Random Scan Gibbs Sampler

As above, let X and Y be discrete r.v.s with joint PMF $p_{X,Y}(x,y)$. We wish to construct a two-dimensional Markov chain (X_n, Y_n) whose stationary distribution is $p_{X,Y}$. Each move of the random scan Gibbs sampler picks a uniformly random component and updates it, according to the conditional distribution given the other component.

Algorithm Random scan Gibbs sampler

Require:

Joint PMF $p_{X,Y}$;

Initial state (X_0, Y_0) ;

Ensure:

Two-dimensional Markov chain (X_n, Y_n) ;

- 1: **repeat**
 - 2: Choose which component to update, with equal probabilities;
 - 3: If the X -component was chosen, draw a value x_{n+1} from the conditional distribution of X given $Y = y_n$, and set $X_{n+1} = x_{n+1}$, $Y_{n+1} = y_n$. Similarly, if the Y -component was chosen, draw a value y_{n+1} from the conditional distribution of Y given $X = x_n$, and set $X_{n+1} = x_n$, $Y_{n+1} = y_{n+1}$;
 - 4: **return** (X_{n+1}, Y_{n+1}) ;
 - 5: **until** $n \geq N$;
-

Random Scan Gibbs as Metropolis-Hastings

Hw

Theorem

The random scan Gibbs sampler is a special case of the Metropolis-Hastings algorithm, in which the proposal is always accepted. In particular, it follows that the stationary distribution of the random scan Gibbs sampler is as desired.

Gibbs Sampling vs. Metropolis-Hastings

- Gibbs sampling emphasizes conditional distributions.
- Metropolis-Hastings emphasizes acceptance probabilities.

Gibbs Sampler for Continuous State Space

In the Gibbs sampler, the target distribution π is an m -dimensional joint density

$$\pi(\mathbf{x}) = \pi(x_1, \dots, x_m).$$

A multivariate Markov chain is constructed whose limiting distribution is π , and which takes values in an m -dimensional space. The algorithm generates elements by iteratively updating each component of an m -dimensional vector conditional on the other $m - 1$ components.

Example: Bivariate Standard Normal Distribution

$$P = \pm 1,$$

Joint PDF $f(x, y) = \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-\frac{x^2 - 2\rho xy + y^2}{2(1-\rho^2)}},$
 $-\infty < x, y < \infty, \quad -1 < \rho < 1,$

$$f(x|Y=y) \sim N(\rho y, 1-\rho^2)$$

$$\xrightarrow{\hspace{1cm}} f(Y|X=x) \sim \underline{N(\rho x, 1-\rho^2)}$$

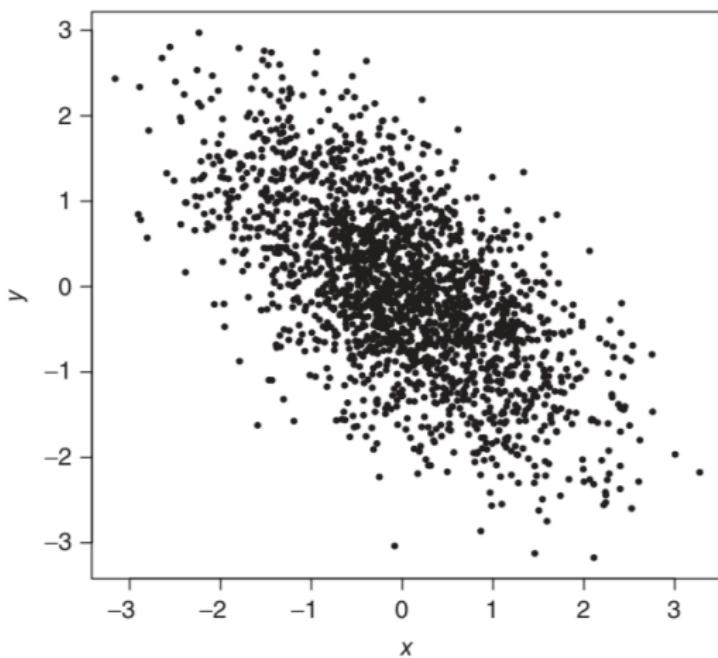
Example: Bivariate Standard Normal Distribution

The Gibbs sampler is implemented to simulate (X, Y) from a bivariate standard normal distribution with correlation ρ .

- ① Initialize: $(x_0, y_0) \leftarrow (0, 0)$, $m \leftarrow 1$.
- ② Generate x_m from the conditional distribution of X given $Y = y_{m-1}$. That is, simulate from a normal distribution with mean ρy_{m-1} and variance $1 - \rho^2$. $f(x|Y=y_{m-1}) \sim N(\rho y_{m-1}, 1-\rho^2)$
- ③ Generate y_m from the conditional distribution of Y given $X = x_m$. That is, simulate from a normal distribution with mean ρx_m and variance $1 - \rho^2$. $f(y|X=x_m) \sim N(\rho x_m, 1-\rho^2)$
- ④ $m \leftarrow m + 1$.
- ⑤ Return to Step 2.

$$\cancel{\sim N(\rho x_{m-1}, 1-\rho^2)}$$

Simulation Results with MCMC



2000

,

$$\rho = -0.6$$

Example: Chicken-Egg with Unknown Parameters

$$N = n, \quad X = x,$$

$$P ? \sim \text{Beta}(a+x, b+n-x)$$

A chicken lays N eggs, where $N \sim \text{Pois}(\lambda)$. Each egg hatches with probability p , where p is unknown; we let $p \sim \text{Beta}(a, b)$. The constants λ, a, b are known.

Here's the catch: we don't get to observe N . Instead, we only observe the number of eggs that hatch, X . Describe how to use Gibbs sampling to find $E(p|X=x)$, the posterior mean of p after observing x hatched eggs.

$$P = p, \quad X = x$$

$$N ?$$

$$\frac{N - X}{P} \sim \text{Pois}(\lambda(1-p))$$

$$N = X + \text{Pois}(\lambda(1-p))$$

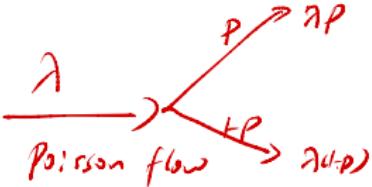
Solution

1^o. # of eggs $N \sim \text{Pois}(\lambda)$;

of hatched eggs.

$$X | P=p \sim \text{Pois}(\lambda p).$$

$$N-X | P=p \sim \text{Pois}(\lambda(1-p)).$$



$$\begin{aligned} f(p | X=x) &\propto \text{Prob.}(X=x | P=p) \cdot \underbrace{f(p)}_{\text{likelihood}} \cdot \underbrace{\text{prior}}_{\propto e^{-\lambda p} \cdot (\lambda p)^x \cdot p^{a-1} (1-p)^{b-1}}. \end{aligned}$$

We can adopt M-H algorithm to generate samples.

2^o. $f_p(p | X=x)$ hard to sample.

$f_{p|N}(p, n | X=x)$ easy to sample.

2.1^o. $f(p | N=n, X=x) \sim \text{Beta}(x+a, n-x+b)$

$f(n | P=p, X=x)$ \sim $X + \text{Pois}(\lambda(1-p))$ shifted poisson

Solution

Solution

We make an initial guess for p and N , then iterate the following steps:

1. Conditional on $N = n$ and $X = x$, draw a new guess for p from the Beta($x + a, n - x + b$) distribution.
2. Conditional on p and $X = x$, the number of unhatched eggs is $Y \sim \text{Pois}(\lambda(1 - p))$ by the chicken-egg story, so we can draw Y from the Pois($\lambda(1 - p)$) distribution and set the new guess for N to be $N = x + Y$.

Solution

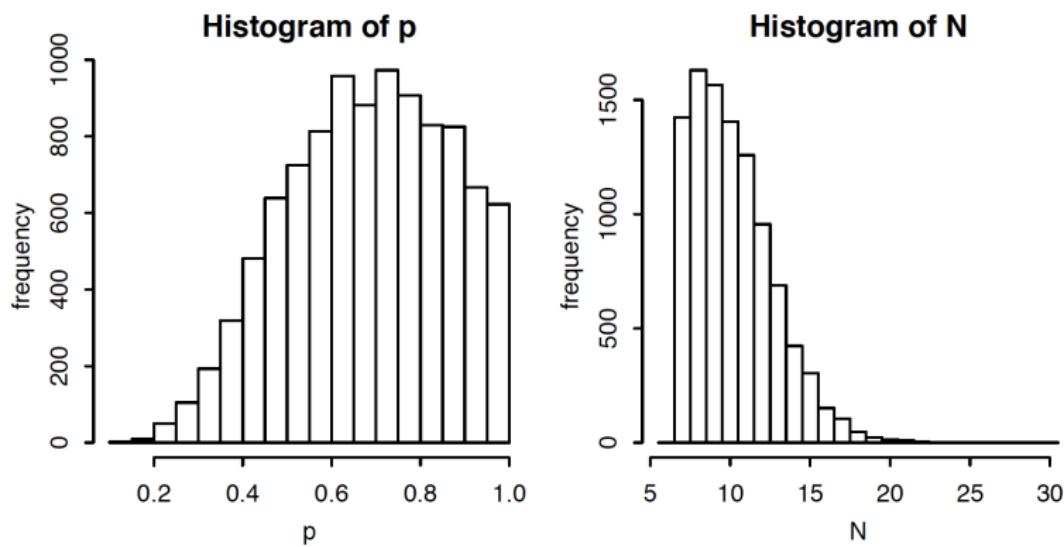


FIGURE 12.5

Histograms of 10^4 draws from the posterior distributions of p and N , where $\lambda = 10$, $a = 1$, $b = 1$, and we observe $X = 7$.

0.68
Sample mean

Example: Three-dimensional Joint Distribution

① PMF $f(x|N=n, p=p) \propto \frac{\binom{n}{x} p^x (1-p)^{n-x}}{n!}$ $\sim \text{Bin}(n, p)$

② PDF $f(p|x=x, N=n) \propto p^x (1-p)^{n-x}$ $\sim \text{Beta}(x+1, n-x+1)$

Random variables X, P and N have joint density

$$\pi(x, p, n) \propto \frac{\binom{n}{x} p^x (1-p)^{n-x} 4^n}{n!}$$

for $x = 0, 1, \dots, n$, $0 < p < 1$, $n = 0, 1, \dots$. The p variable is continuous; x and n are discrete.

③ PMF $f(n|x=x, p=p) \propto \frac{\binom{n}{x} (1-p)^{n-x} 4^n}{n!}$

Shifted Poisson. $= x + z$

$z \sim \text{Pois}(4(1-p))$

Solution

Solution

The Gibbs sampler, with arbitrary initial value, is implemented as follows:

1. Initialize: $(x_0, p_0, n_0) \leftarrow (1, 0.5, 2)$

Systematic Scan

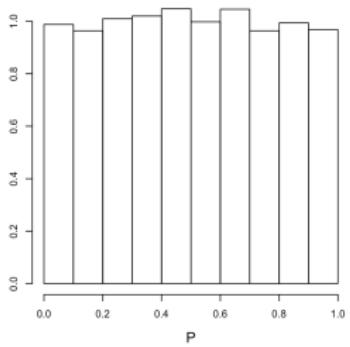
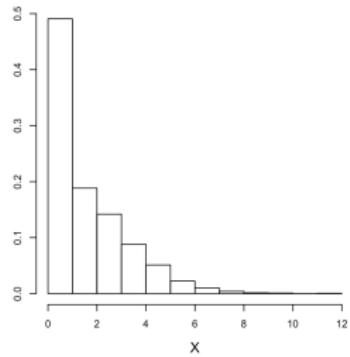
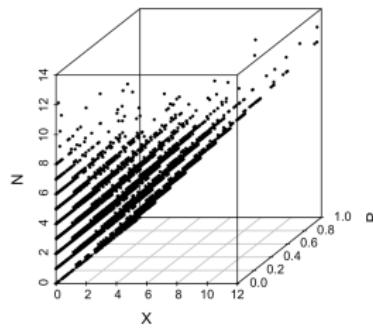
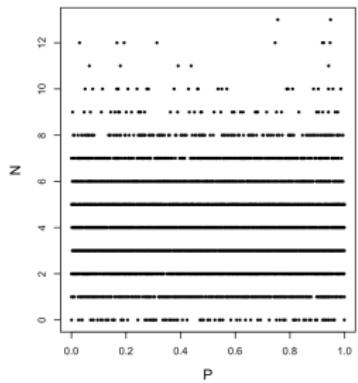
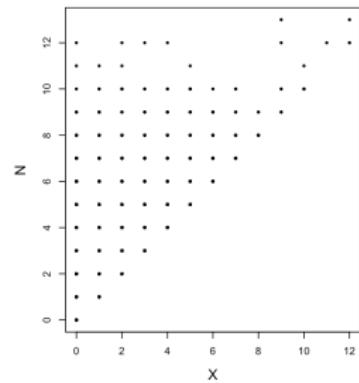
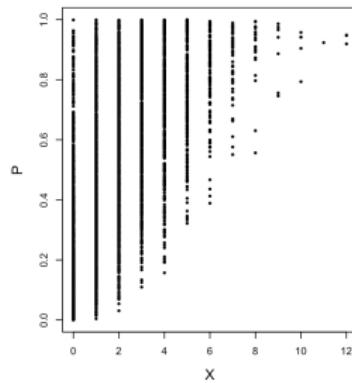
$$m \leftarrow 1$$

2. Generate x_m from a binomial distribution with parameters n_{m-1} and p_{m-1}
3. Generate p_m from a beta distribution with parameters $x_m + 1$ and $n_{m-1} - x_m + 1$.
4. Let $n_m = z + \epsilon_m$ where z is simulated from a Poisson distribution with parameter $4(1 - p_m)$.
5. $m \leftarrow m + 1$
6. Return to Step 2.

The output of the Gibbs sampler is a sequence of samples

$$(X_0, P_0, N_0), (X_1, P_1, N_1), (X_2, P_2, N_2), \dots$$

Simulation Results with MCMC



Outline

- 1 Introduction of MCMC
- 2 Metropolis–Hastings Algorithm
- 3 Gibbs Sampler
- 4 Hamiltonian Monte Carlo Method
- 5 References

Continuous
PDF

Motivation: Random-Walk Metropolis

Let $Y|\theta \sim \mathcal{N}(\theta, \sigma^2)$, where σ^2 is known but θ is unknown. Using the Bayesian framework, we treat θ as a random variable, with prior given by $\theta \sim \mathcal{N}(\mu, \tau^2)$ for some known constants μ and τ^2 . That is, we have the two-level model

$$\begin{aligned}\theta &\sim \mathcal{N}(\mu, \tau^2) \\ Y|\theta &\sim \mathcal{N}(\theta, \sigma^2) \\ \theta &= \theta\end{aligned}$$

Describe how to use the Metropolis-Hastings algorithm to find the posterior mean and variance of θ after observing the value of Y .

Solution P. $f_{\theta|Y=y}(\theta|y) \propto f_{Y|\theta}(y|\theta) \cdot f_{\theta}(\theta)$

$$\propto e^{-\frac{1}{2\sigma^2}(y-\theta)^2} \cdot e^{-\frac{1}{2\tau^2}(\theta-\mu)^2}$$

Because of Normal-Normal Conjugacy,

$$\theta|Y=y \sim N\left(\underbrace{\frac{\frac{1}{\sigma^2}}{\frac{1}{\sigma^2} + \frac{1}{\tau^2}} y + \frac{\frac{1}{\tau^2}}{\frac{1}{\sigma^2} + \frac{1}{\tau^2}} \mu}_{\text{Posterior mean of } \theta|Y=y}, \underbrace{\frac{1}{\frac{1}{\sigma^2} + \frac{1}{\tau^2}}}_{\text{Posterior Variance of } \theta|Y=y}\right)$$

Posterior mean of
 $\theta|Y=y$

Posterior Variance
of
 $\theta|Y=y$

Solution 2^o. M-H.

$$\propto e^{-\frac{1}{2\sigma^2}(y-x)^2} \cdot e^{-\frac{1}{2t^2}(t\theta H)^2}$$

construct M.C. whose stationary distribution $(f_{\theta|Y}(y))$

Generates $\theta_0, \theta_1, \dots$

(a) if $\theta_n = x$, propose a new state x' [rw]

$$x' = x + \varepsilon_n \quad \varepsilon_n \sim N(0, d^2)$$

$$= x + \underline{\text{d}}. \underline{N(0, 1)}$$

$$\Rightarrow x = x' - d.N(0, 1)$$

$$\stackrel{D}{=} x' + d.N(0, 1)$$

(b). the acceptance prob.

$$\alpha(x, x') = \min \left(\frac{x' f(x, x)}{x f(x', x)}, 1 \right)$$

$$= \min \left(\frac{\cancel{f_{\theta|Y}(x'|y)}}{\cancel{f_{\theta|Y}(x|y)}}, 1 \right)$$

$$\Rightarrow f(x, x')$$

$$\stackrel{D}{=} f(x', x)$$

Solution

Solution

Simulation Results with MH

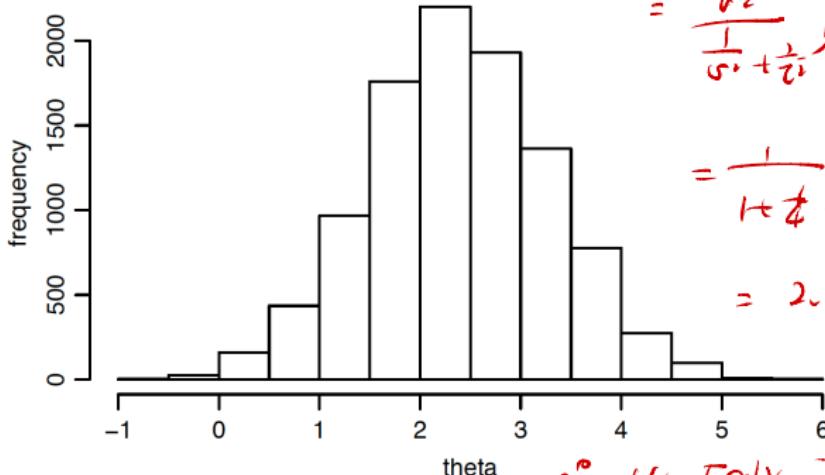
Theory

$$1^{\circ}. E[\theta | Y=y]$$

$$= \frac{\frac{1}{\sigma^2}}{\frac{1}{\sigma^2} + \frac{1}{\tau^2}} y + \frac{\frac{1}{\tau^2}}{\frac{1}{\sigma^2} + \frac{1}{\tau^2}} \mu$$

$$= \frac{1}{1+\frac{1}{4}} \cdot 3 + \frac{\frac{1}{4}}{1+\frac{1}{4}} \cdot 0$$

$$= 2.4;$$



$$2^{\circ}. \text{Var}[\theta | Y=y] = \frac{1}{\frac{1}{\sigma^2} + \frac{1}{\tau^2}}$$

$$= \frac{1}{1+\frac{1}{4}} = 0.8$$

FIGURE 12.2

Histogram of 10^4 draws from the posterior distribution of θ given $Y=3$, obtained using Metropolis-Hastings with $\mu=0$, $\sigma^2=1$, and $\tau^2=4$. The sample mean is 2.4 and the sample variance is 0.8, in agreement with the theoretical values.

Simulation Results with MH

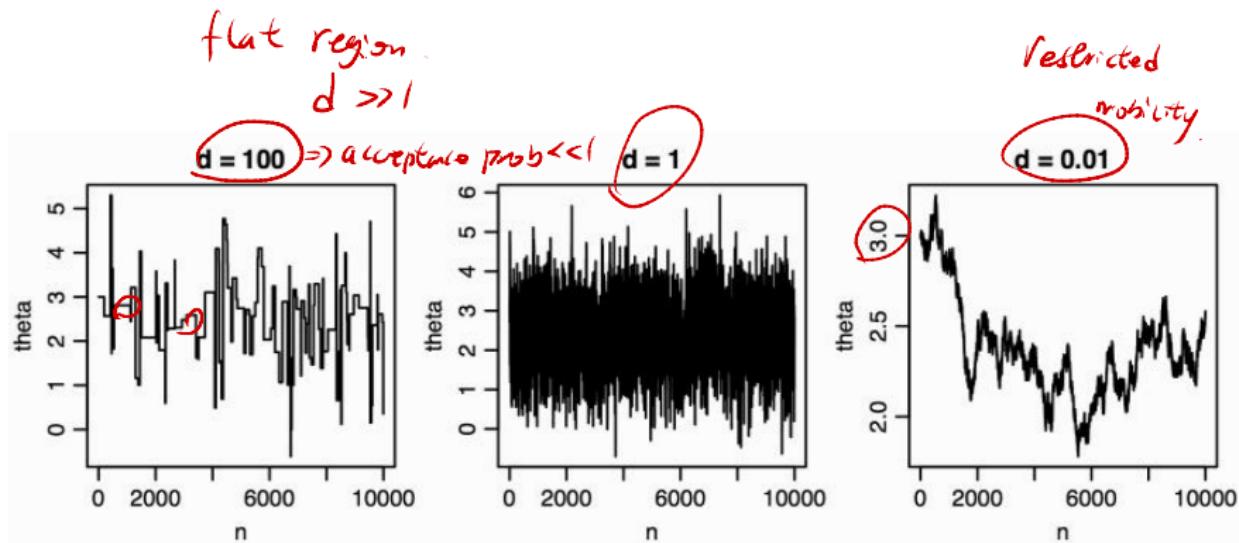


FIGURE 12.3

Trace plots of θ_n as a function of the iteration number n , for $d = 100, 1, 0.01$.

Random Walk Metropolis: RWM

- Continuous state space with target distribution π

- Proposal state:

$$\hat{X}_{n+1} = X_n + \eta \cdot \mathcal{N}(0, I)$$

(I)

n -dimension
identity matrix

- One-step transition probability density: $f_{x,y} = f_{y,x}$

- Accept the proposal state with probability

$$a_{X_n, \hat{X}_{n+1}} = \min \left\{ 1, \frac{\pi(\hat{X}_{n+1})}{\pi(X_n)} \right\}.$$

- Zero-order distribution information with dumb exploration of state space & low acceptance probability in general

Gradient Information

- Explore the state space with gradient information
- Gradient ascent algorithm to find the local maximum of function g :

$$x_{n+1} = x_n + \eta \nabla g(x_n).$$

- Add stochastic noise to avoid stuck in local maximum:

$$x_{n+1} = x_n + \eta \nabla g(x_n) + \sqrt{2\eta} \cdot \mathcal{N}(0, I).$$

- Why $\sqrt{2}$? : Overdamped Langevin Ito diffusion

$$dX_t = \nabla g(X_t) dt + \sqrt{2} dB_t \rightarrow \text{Brownian Motion}$$

with Euler-Maruyama discrete approximation

$$x_{n+1} = x_n + \eta \nabla g(x_n) + \sqrt{2\eta} \cdot \mathcal{N}(0, I)$$

Gradient Information

- Now the target distribution is π
- Choice of g , π or $\log(\pi)$?
- The First Choice:

$$\max_x \pi(x) \Leftrightarrow \max_x (\log \pi(x))$$

$$S_x \propto e^{\pi(x)}$$

$$x_{n+1} = x_n + \eta \nabla \pi(x_n) + \sqrt{2\eta} \cdot \mathcal{N}(0, I).$$

- The Second Choice:

$$x_{n+1} = x_n + \eta \nabla \log \pi(x_n) + \sqrt{2\eta} \cdot \mathcal{N}(0, I).$$

- Which one is better?

$$\nabla \log \pi(x)$$

$$\pi(x) \propto g(x)$$

$$S_x \propto e^{\frac{\log \pi(x)}{\pi(x)}}$$

$$\pi(x) = \frac{1}{Z} g(x)$$

$$\log \pi(x) = \log Z + \log g(x)$$

$$\nabla \log \pi(x) = \nabla \log g(x)$$

Langevin Monte Carlo: LMC

- Continuous state space with target distribution π
- Proposal state:

$$\hat{X}_{n+1} = X_n + \eta \nabla \log \pi(x_n) + \sqrt{2\eta} \cdot \mathcal{N}(0, I)$$

- One-step transition probability density:

$$f_{x,y} \propto \exp \left(-\frac{1}{4\eta} (y - x - \eta \nabla \log \pi(x))^2 \right).$$

- Accept the proposal state with probability

$$a_{X_n, \hat{X}_{n+1}} = \min \left\{ 1, \frac{\pi(\hat{X}_{n+1}) f_{\hat{X}_{n+1}, X_n}}{\pi(X_n) f_{X_n, \hat{X}_{n+1}}} \right\}.$$

- Also called Metropolis-adjusted Langevin algorithm (MALA)

Langevin Monte Carlo: LMC

- Usually proposes moves into regions of higher density (π): more likely to be accepted
- Optimal acceptance probability is 0.574 for a limited classes of target distributions
- First-order distribution information with smart exploration of state space & mid acceptance probability in general

Hamiltonian Monte Carlo: HMC

- Random Walk Metropolis (RWM): Zero-order distribution information with dumb exploration of state space & low acceptance probability in general
- Langevin Monte Carlo(LMC): First-order distribution information with smart exploration of state space & mid acceptance probability in general
- Can we find MCMC methods with smart exploration of state space & high acceptance probability in general?
- Yes, We Can!

Hamiltonian Dynamics

- Classical mechanics (or Newtonian mechanics):

$$\underline{\underline{F = ma = m \frac{d^2 \mathbf{x}_t}{dt^2}}}.$$

- Lagrangian mechanics with Lagrangian:

$$L(\underline{\mathbf{x}_t}, \underline{\dot{\mathbf{x}}_t}, t) = \underline{V} - \underline{U} = \text{Kinetic Energy} - \text{Potential Energy}.$$

- Tautochrone curve problem motivated Euler-Lagrange equation and the calculus of variations
- Hamiltonian mechanics(also called Hamiltonian dynamics): reformulation of Lagrangian mechanics with momenta

Hamiltonian Dynamics

- d -dimensional space
- x : position vector of the moving object
- $U(x)$: potential energy
- ω : momentum vector of the moving object
- $V(\omega)$: Kinetic Energy
- $H(x, \omega) = U(x) + V(\omega)$: Hamiltonian(energy)
- $\{(x, \omega)\}$: phase space(coordinate system that is defined in terms of position and momentum)
- Hamiltonian equations: $j = 1, \dots, d$

$$\frac{dH(x, \omega)}{dt} = \sum_{j=1}^d \left\{ \frac{\partial H}{\partial x_j} \cdot \frac{dx_j}{dt} + \frac{\partial H}{\partial \omega_j} \cdot \frac{d\omega_j}{dt} \right\}$$

$$= \sum_{j=1}^d \left\{ \frac{\partial H}{\partial x_j} \cdot \frac{\partial H}{\partial \omega_j} - \frac{\partial H}{\partial \omega_j} \cdot \frac{\partial H}{\partial x_j} \right\}$$

相

$$= 0$$

$$\Rightarrow \underline{H(x, \omega)} = \underline{\text{const}}$$

$$\begin{aligned}\frac{dx_j}{dt} &= \frac{\partial H}{\partial \omega_j} \\ \frac{d\omega_j}{dt} &= -\frac{\partial H}{\partial x_j}.\end{aligned}$$

Hamiltonian Dynamics: Numerical Solution

- No analytical solution for Hamiltonian equations
- Approximation with Taylor expansion: $\forall j = 1, \dots, d$

0, 8, 28, 36

$$x_j(t + \delta) \approx x_j(t) + \delta \frac{dx_j}{dt} = x_j(t) + \delta \frac{\partial H(\mathbf{x}_t, \boldsymbol{\omega}_t)}{\partial \omega_j} = x_j(t) + \delta \frac{\partial V(\boldsymbol{\omega}_t)}{\partial \omega_j}$$

$$\omega_j(t + \delta) \approx \omega_j(t) + \delta \frac{d\omega_j}{dt} = \omega_j(t) - \delta \frac{\partial H(\mathbf{x}_t, \boldsymbol{\omega}_t)}{\partial x_j} = \omega_j(t) - \delta \frac{\partial U(\mathbf{x}_t)}{\partial x_j}$$

- Three methods: Euler's Method, Modified Euler's Method, The Leapfrog Method

Numerical Solution I: Euler's Method

x_j, ω_j both updated.

$$\omega_j(t + \delta) = \omega_j(t) - \delta \frac{\partial U(x_t)}{\partial x_j}$$

instability

$$x_j(t + \delta) = x_j(t) + \delta \frac{\partial V(\omega_t)}{\partial \omega_j}.$$

Numerical Solution II: Modified Euler's Method

$$\begin{aligned}\omega_j(t + \delta) &= \omega_j(t) - \delta \frac{\partial U(\underline{x}_t)}{\partial x_j} \\ \underline{x}_j(t + \delta) &= x_j(t) + \delta \frac{\partial V(\underline{\omega}_{t+\delta})}{\partial \omega_j}.\end{aligned}$$

Alternate updated

stable

Numerical Solution III: The Leapfrog Method

$$\underline{\omega_j \left(t + \frac{\delta}{2} \right)} = \omega_j(t) - \frac{\delta}{2} \cdot \frac{\partial U(\mathbf{x}_t)}{\partial x_j}$$

half step

$$x_j(t + \delta) = x_j(t) + \delta \cdot \frac{\partial V(\omega_{t+\frac{\delta}{2}})}{\partial \omega_j}$$

one-step

$$\underline{\omega_j(t + \delta)} = \omega_j \left(t + \frac{\delta}{2} \right) - \frac{\delta}{2} \cdot \frac{\partial U(\mathbf{x}_{t+\delta})}{\partial x_j}$$

half step.

Example

- $d = \underline{1}$ initial position $x_0 = \underline{0}$, initial momentum $\omega_0 = \underline{1}$
- $H(x, \omega) = U(x) + V(\omega) = \frac{x^2}{2} + \frac{\omega^2}{2}$
- Phase space:

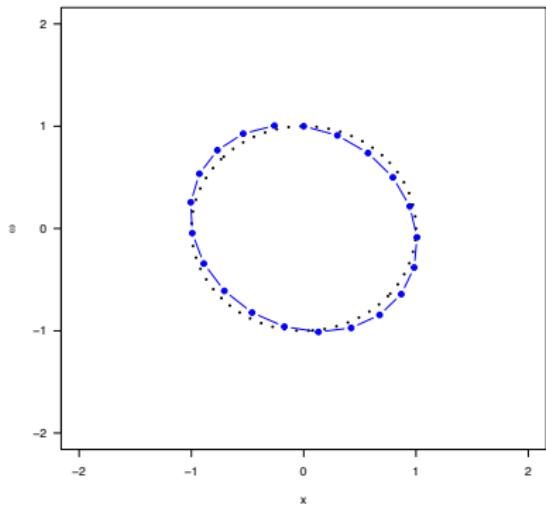
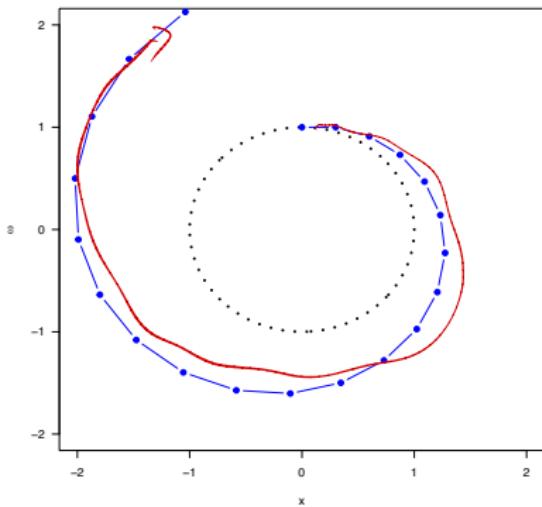
$$H(x_t, \omega_t) = \frac{x_t^2}{2} + \frac{\omega_t^2}{2} = H(x_0, \omega_0) = \frac{x_0^2}{2} + \frac{\omega_0^2}{2} = \frac{1}{2}.$$

- Equivalent expression:

$$x_t^2 + \omega_t^2 = 1, \forall t \geq 0$$

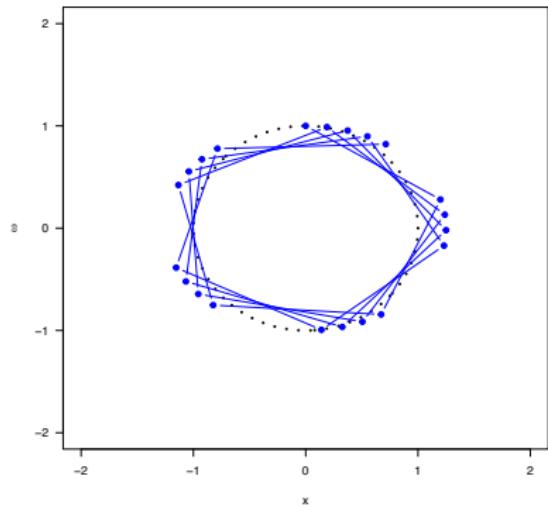
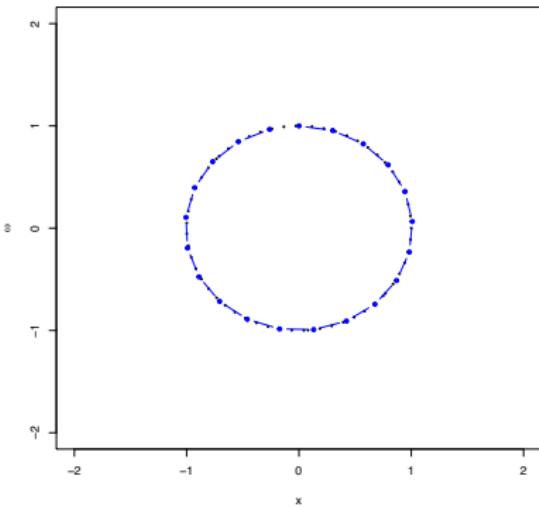
Example

- Step size $\delta = 0.3$ with 20 steps
- LHS is the result with Euler's Method
- RHS is the result with Modified Euler's Method



Example

- Both adopt The Leapfrog Method
- LHS is the result with step size $\delta = 0.3$ & 20 steps
- RHS is the result with step size $\delta = 1.2$ & 20 steps



Hamiltonian Monte Carlo: HMC

- Main variable: position $x \in \mathbb{R}^d$ of the moving object
- Target distribution $\pi(x)$ is the distribution of the position
- Auxiliary variable: momentum $\omega \in \mathbb{R}^d$ of the moving object
- State of the Markov process: (x, ω)
- State space: phase space of corresponding Hamiltonian Dynamics
- Stationary distribution (joint pdf) over the state space:

$$p(x, \omega) \propto e^{-H(x, \omega)} = e^{-U(x) - V(\omega)} = e^{-U(x)} \cdot e^{-V(\omega)}.$$

Hamiltonian Monte Carlo: HMC

- Stationary distribution (joint pdf) over the state space:

$$p(\mathbf{x}, \boldsymbol{\omega}) \propto e^{-H(\mathbf{x}, \boldsymbol{\omega})} = e^{-U(\mathbf{x}) - V(\boldsymbol{\omega})} = \underbrace{e^{-U(\mathbf{x})}}_{\pi(\mathbf{x})} \cdot \underbrace{e^{-V(\boldsymbol{\omega})}}_{\pi(\boldsymbol{\omega})} = \pi(\mathbf{x}) \cdot \pi(\boldsymbol{\omega}).$$

- Set $e^{-U(\mathbf{x})} = \pi(\mathbf{x})$, then

$$\underbrace{U(\mathbf{x})}_{- \log(\pi(\mathbf{x}))} = - \log(\pi(\mathbf{x})).$$

- Set $e^{-V(\boldsymbol{\omega})} = \pi(\boldsymbol{\omega})$, usually choose multivariate Normal distribution:

$$\pi(\boldsymbol{\omega}) \propto e^{-\frac{1}{2}\boldsymbol{\omega}^T \Sigma^{-1} \boldsymbol{\omega}} \sim \mathcal{N}(0, \Sigma)$$

where

$$\Sigma = \underbrace{\text{diag}(m_1, \dots, m_d)}_{\text{mass matrix}}$$

then

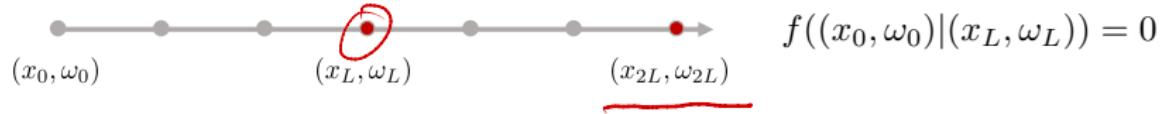
$$V(\boldsymbol{\omega}) = \sum_{j=1}^d \frac{\omega_j^2}{2m_j}, \forall m_j > 0.$$

One-Round Hamiltonian Monte Carlo

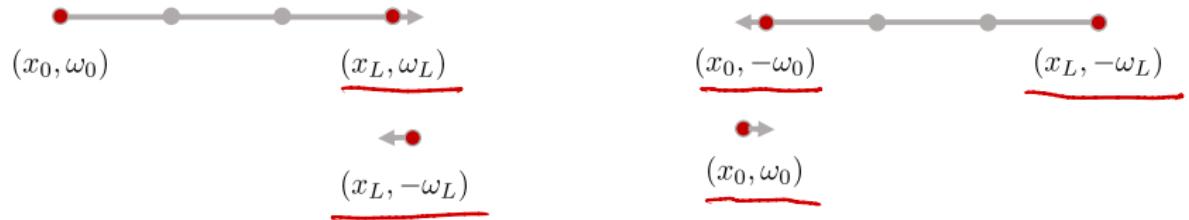
- Start with a given initial position x_0
- Set the initial momentum ω_0 as one sample from $\mathcal{N}(0, \Sigma)$
- Apply the leapfrog method with initial state (x_0, ω_0) with step-size δ and step number L ; then record state (x_L, ω_L)
- Reverse the direction of momentum and obtain the state $(x_L, -\omega_L)$.
- Set the state $(x_L, -\omega_L)$ as the proposal state.
- Accept the proposal state with probability

$$\begin{aligned} a_{(x_0, \omega_0), (x_L, -\omega_L)} &= \min \left\{ 1, \frac{\exp(-H(x_L, -\omega_L))}{\exp(-H(x_0, \omega_0))} \right\} \\ &= \min \left\{ 1, \frac{\exp(-U(x_L) - V(-\omega_L))}{\exp(-U(x_0) - V(\omega_0))} \right\}. \end{aligned}$$

Why Reverse the Direction of Momentum?



Why Reverse the Direction of Momentum?



$$f((x_L, -\omega_L) | (x_0, \omega_0)) = f((x_0, \omega_0) | (x_L, -\omega_L)) > 0$$

Acceptance Probability is Close to 1

In Theory

$$H(x_L, -\omega_L) = H(x_0, \omega_0)$$

discrete
approximation

$$H(x_L, -\omega_L) \approx H(x_0, \omega_0)$$

$$\begin{aligned} a_{(x_0, \omega_0), (x_L, -\omega_L)} &= \min \left\{ 1, \frac{p(x_L, -\omega_L) f((x_0, \omega_0) | (x_L, -\omega_L))}{p(x_0, \omega_0) f((x_L, -\omega_L) | (x_0, \omega_0))} \right\} \\ &= \min \left\{ 1, \frac{p(x_L, -\omega_L)}{p(x_0, \omega_0)} \right\} \\ &= \min \left\{ 1, \frac{\exp(-H(x_L, -\omega_L))}{\exp(-H(x_0, \omega_0))} \right\} \\ &= \min \left\{ 1, \frac{\exp(-U(x_L) - V(-\omega_L))}{\exp(-U(x_0) - V(\omega_0))} \right\}. \end{aligned}$$

≈ 1

Example I

$g(x)$

$$f(x) \propto \underbrace{(1 + \frac{1}{5}x^2)^{-3}}$$

$$U(x) = -\log g(x)$$

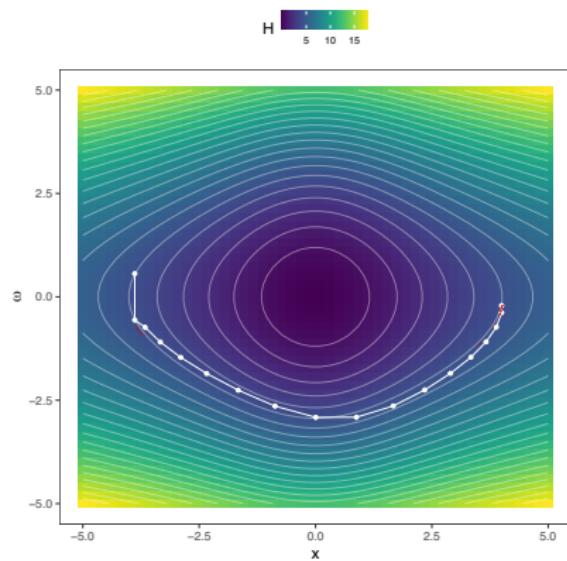
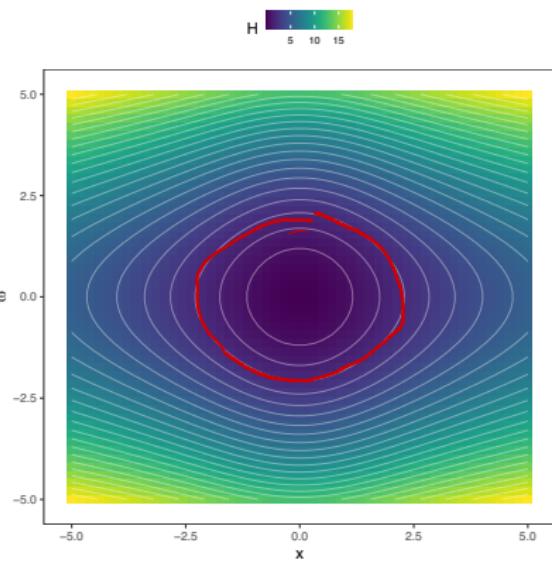
- Target distribution(PDF):

$$\underline{f(x) = \frac{8}{3\pi\sqrt{5}} \left(1 + \frac{1}{5}x^2\right)^{-3}, -\infty < x < \infty.} \quad = 3 \log(1 + \frac{1}{5}x^2)$$

$$V(\omega) = \frac{1}{2}\omega^2, \quad \underline{\omega \sim N(0, 1)}$$

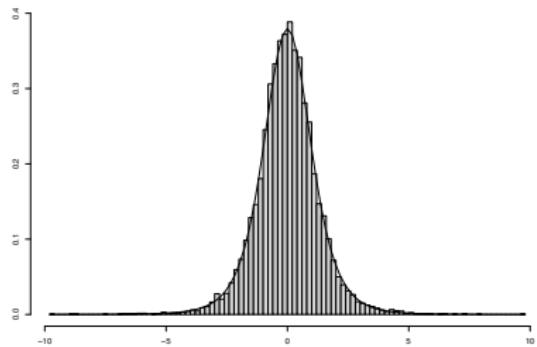
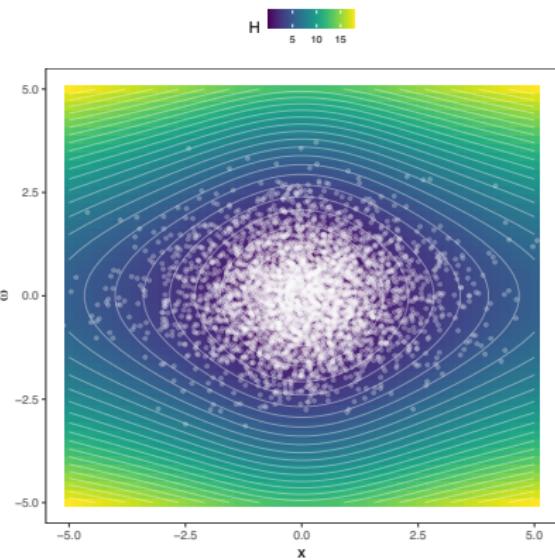
Example I

$$\delta=0.3, \quad L=15$$



Example I

$T = 10000$



Example II: Gaussian Mixture Distribution

$$f(x) \propto 0.4 e^{-\frac{1}{2}(x-2)^2} + 0.6 e^{-\frac{1}{2}(x+2)^2}$$

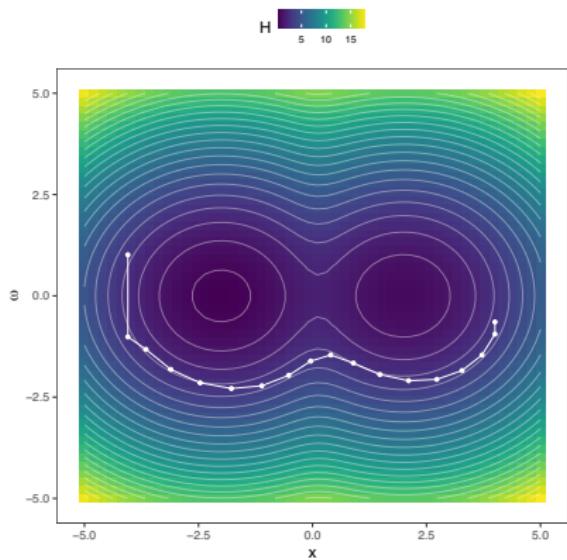
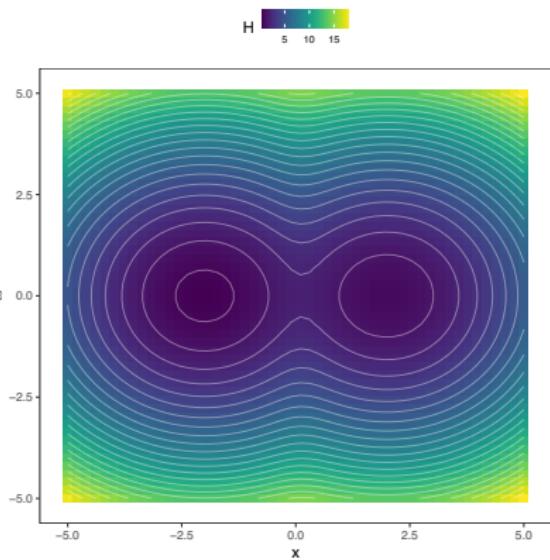
$\underbrace{\hspace{10em}}_{g(x)}$

- Target distribution(PDF): $U(x) = -(\log g(x)) = -\log(g(x))$

$$f(x) = 0.4 \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-2)^2} + 0.6 \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x+2)^2}, \quad -\infty < x < \infty.$$

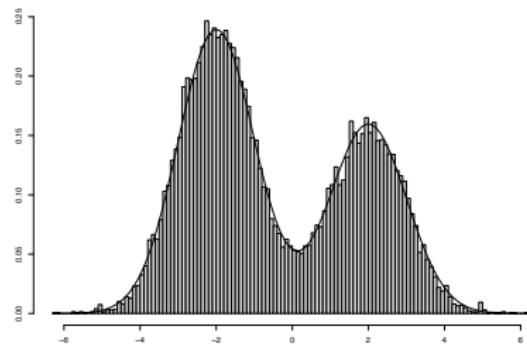
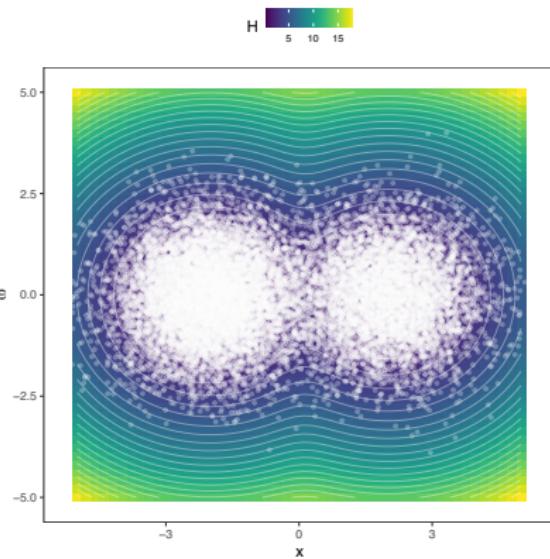
$$U(\omega) = \frac{1}{2}\omega^2, \quad \underline{\omega \in N(0,1)}$$

Example II



Example II

$T = 20000$

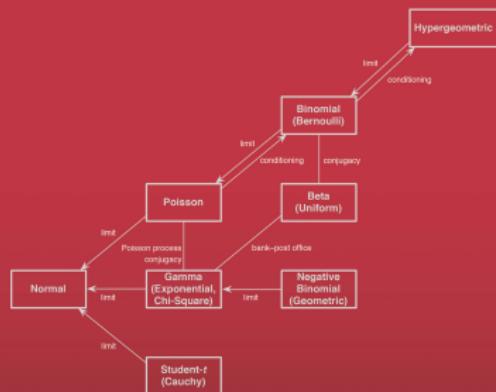


Outline

- 1 Introduction of MCMC
- 2 Metropolis–Hastings Algorithm
- 3 Gibbs Sampler
- 4 Hamiltonian Monte Carlo Method
- 5 References

Texts in Statistical Science

Introduction to Probability



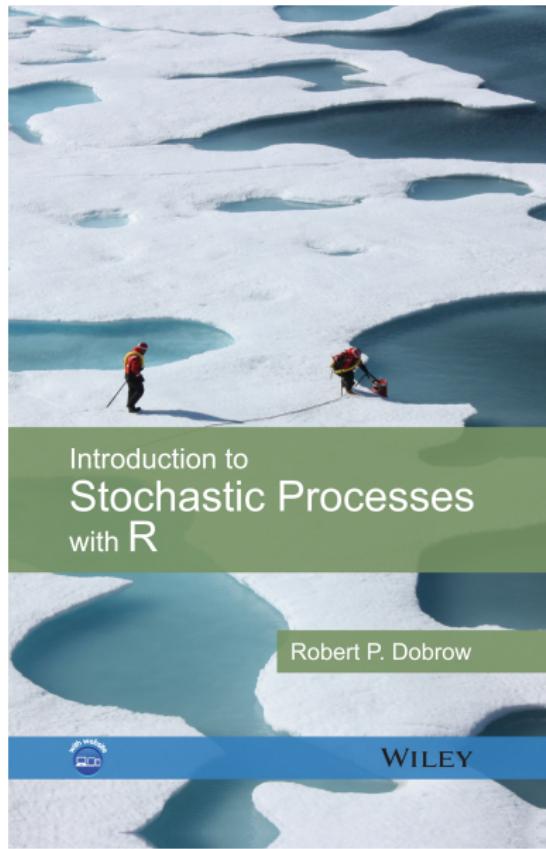
Joseph K. Blitzstein
Jessica Hwang



CRC Press
Taylor & Francis Group
A CHAPMAN & HALL BOOK

BH

- Introduction to Probability
- Chapman & Hall/CRC, 2014.
- Chapman & Hall/CRC, 2019.
- Chapter 12



- SPR
- Introduction to Stochastic Processes with R
 - John Wiley & Son, 2016.
 - Chapter 5

Use R!

Christian P. Robert
George Casella

Introducing Monte Carlo Methods with R

 Springer

RC

- Introducing Monte Carlo Methods with R
- Springer, 2010.
- All chapters including Chapter 8 (when to stop MCMC algorithm)

Others

SOTA

- Betancourt, Michael (2018). “A Conceptual Introduction to Hamiltonian Monte Carlo”. arXiv:1701.02434
- Neal, Radford M (2011). “MCMC Using Hamiltonian Dynamics”