

LECTURE 16: MARKOV CHAIN MONTE CARLO (CONTD)

STAT 545: INTRO. TO COMPUTATIONAL STATISTICS

Vinayak Rao

Purdue University

November 10, 2016

MARKOV CHAIN MONTE CARLO

We are interested in a distribution $\pi(x) = \frac{f(x)}{Z}$

(e.g. want the mean, quantiles etc.)

Monte Carlo: approximate with independent samples from π

MCMC: produce dependent samples via a Markov chain

$$X_0 \rightarrow X_1 \rightarrow X_2 \rightarrow X_3 \rightarrow \cdots \rightarrow X_{N-1} \rightarrow X_N$$

Use dependent samples to approximate integrals w.r.t. $\pi(x)$:

$$\frac{1}{N} \sum_{i=1}^N g(x_i) \approx \mathbb{E}_{\pi}[g] \quad \text{as}$$

Let $\mathcal{T}(x_i \rightarrow x_{i+1})$ be the Markov transition kernel. We require:

Let $\mathcal{T}(x_i \rightarrow x_{i+1})$ be the Markov transition kernel. We require:

- π is a stationary distribution of \mathcal{T} :

$$\pi(x) = \int_{\mathcal{X}} \pi(y) \mathcal{T}(y \rightarrow x) dy$$

Let $\mathcal{T}(x_i \rightarrow x_{i+1})$ be the Markov transition kernel. We require:

- π is a stationary distribution of \mathcal{T} :

$$\pi(x) = \int_{\mathcal{X}} \pi(y) \mathcal{T}(y \rightarrow x) dy$$

Also require that \mathcal{T} be

- irreducible: not stuck to some part of \mathcal{X} forever

Let $\mathcal{T}(x_i \rightarrow x_{i+1})$ be the Markov transition kernel. We require:

- π is a stationary distribution of \mathcal{T} :

$$\pi(x) = \int_{\mathcal{X}} \pi(y) \mathcal{T}(y \rightarrow x) dy$$

Also require that \mathcal{T} be

- irreducible: not stuck to some part of \mathcal{X} forever
- aperiodic: not stuck to some part of \mathcal{X} for e.g. even iterations. Can fix this with a 'lazy' Markov chain that allows self-transitions.

Let $\mathcal{T}(x_i \rightarrow x_{i+1})$ be the Markov transition kernel. We require:

- π is a stationary distribution of \mathcal{T} :

$$\pi(x) = \int_{\mathcal{X}} \pi(y) \mathcal{T}(y \rightarrow x) dy$$

Also require that \mathcal{T} be

- irreducible: not stuck to some part of \mathcal{X} forever
- aperiodic: not stuck to some part of \mathcal{X} for e.g. even iterations. Can fix this with a 'lazy' Markov chain that allows self-transitions.

Finally, for infinite state-spaces (e.g. the real line), need an additional condition:

- positive recurrent: revisits every neighborhood infinitely often

With these conditions, our chain is *ergodic*

For any initialization:

$$\frac{1}{N} \sum_{i=1}^N g(x_i) \rightarrow \mathbb{E}_{\pi}[g] \quad \text{as } N \rightarrow \infty \quad (\text{Ergodicity})$$

With these conditions, our chain is *ergodic*

For any initialization:

$$\frac{1}{N} \sum_{i=1}^N g(x_i) \rightarrow \mathbb{E}_{\pi}[g] \quad \text{as } N \rightarrow \infty \quad (\text{Ergodicity})$$

We eventually forget the arbitrary initialization.

Typically, we discard the first B burn-in samples.

With these conditions, our chain is *ergodic*

For any initialization:

$$\frac{1}{N} \sum_{i=1}^N g(x_i) \rightarrow \mathbb{E}_{\pi}[g] \quad \text{as } N \rightarrow \infty \quad (\text{Ergodicity})$$

We eventually forget the arbitrary initialization.

Typically, we discard the first B burn-in samples.

A good transition kernel has:

- A short burn-in period.
- Fast mixing (small dependence across samples).

The Markov transition kernel \mathcal{T} must satisfy

$$\pi(x_{n+1}) = \int_{\mathcal{X}} \pi(x_n) \mathcal{T}(x_{n+1}|x_n) dx_n$$

The Markov transition kernel \mathcal{T} must satisfy

$$\pi(x_{n+1}) = \int_{\mathcal{X}} \pi(x_n) \mathcal{T}(x_{n+1}|x_n) dx_n$$

Usually, we enforce the stronger condition of detailed balance:

$$\pi(x_{n+1}) \mathcal{T}(x_n|x_{n+1}) = \pi(x_n) \mathcal{T}(x_{n+1}|x_n)$$

(Sufficient but not necessary)

THE PROBLEM

Given some probability density $\pi(x) = f(x)/Z$:

- How do you construct a transition kernel \mathcal{T} with π as it's stationary distribution?
- How do you construct a *good* transition kernel

Focus of a huge literature.

THE PROBLEM

Given some probability density $\pi(x) = f(x)/Z$:

- How do you construct a transition kernel \mathcal{T} with π as it's stationary distribution?
- How do you construct a *good* transition kernel

Focus of a huge literature.

One approach: the Metropolis-Hastings algorithm

THE METROPOLIS-HASTINGS ALGORITHM

The simplest and most widely applicable MCMC algorithm. Featured in Dongarra & Sullivan (2000)'s list of top 10 algorithms.

1. Metropolis Algorithm for Monte Carlo
2. Simplex Method for Linear Programming
3. Krylov Subspace Iteration Methods
4. The Decompositional Approach to Matrix Computations
5. The Fortran Optimizing Compiler
6. QR Algorithm for Computing Eigenvalues
7. Quicksort Algorithm for Sorting
8. Fast Fourier Transform
9. Integer Relation Detection
10. Fast Multipole Method

THE METROPOLIS-HASTINGS ALGORITHM

A random walk algorithm

Choose a proposal distrib. $q(x_{new}|x_{old})$. E.g. $x_{new} \sim \mathcal{N}(x_{old}, \sigma^2 I)$

THE METROPOLIS-HASTINGS ALGORITHM

A random walk algorithm

Choose a proposal distrib. $q(x_{new}|x_{old})$. E.g. $x_{new} \sim \mathcal{N}(x_{old}, \sigma^2 I)$

Initialize chain at some starting point x_0 .

Repeat:

- Propose a new point x^* according to $q(x^*|x_n)$.
- Define $\alpha = \min \left(1, \frac{\pi(x^*)q(x_n|x^*)}{\pi(x_n)q(x^*|x_n)} \right) = \min \left(1, \frac{f(x^*)q(x_n|x^*)}{f(x_n)q(x^*|x_n)} \right)$
- Set $x_{n+1} = x^*$ with probability α , else $x_{n+1} = x_n$.

THE METROPOLIS-HASTINGS ALGORITHM

A random walk algorithm

Choose a proposal distrib. $q(x_{new}|x_{old})$. E.g. $x_{new} \sim \mathcal{N}(x_{old}, \sigma^2 I)$

Initialize chain at some starting point x_0 .

Repeat:

- Propose a new point x^* according to $q(x^*|x_n)$.
- Define $\alpha = \min \left(1, \frac{\pi(x^*)q(x_n|x^*)}{\pi(x_n)q(x^*|x_n)} \right) = \min \left(1, \frac{f(x^*)q(x_n|x^*)}{f(x_n)q(x^*|x_n)} \right)$
- Set $x_{n+1} = x^*$ with probability α , else $x_{n+1} = x_n$.

Comments:

- Do not need to calculate the normalization constant Z .
- Accept/reject steps ensure this has the correct distribution.

THE METROPOLIS-HASTINGS ALGORITHM

A random walk algorithm

Choose a proposal distrib. $q(x_{new}|x_{old})$. E.g. $x_{new} \sim \mathcal{N}(x_{old}, \sigma^2 I)$

Initialize chain at some starting point x_0 .

Repeat:

- Propose a new point x^* according to $q(x^*|x_n)$.
- Define $\alpha = \min \left(1, \frac{\pi(x^*)q(x_n|x^*)}{\pi(x_n)q(x^*|x_n)} \right) = \min \left(1, \frac{f(x^*)q(x_n|x^*)}{f(x_n)q(x^*|x_n)} \right)$
- Set $x_{n+1} = x^*$ with probability α , else $x_{n+1} = x_n$.

Comments:

- Do not need to calculate the normalization constant Z .
- Accept/reject steps ensure this has the correct distribution.
- On rejection, keep old sample (i.e. there will be repetition)

THE METROPOLIS-HASTINGS ALGORITHM

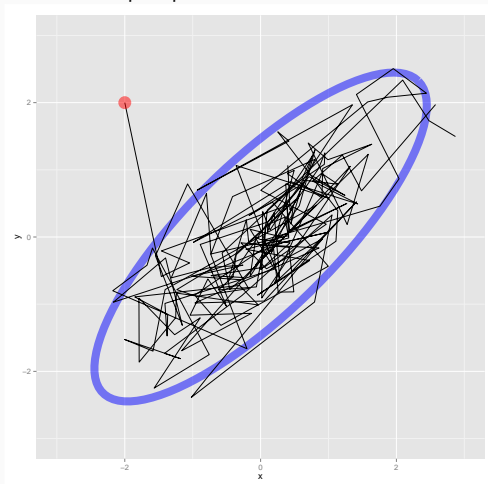
For a symmetric proposal ($q(x^*|x_n) = q(x_n|x^*)$):

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

The Metropolis algorithm.

THE METROPOLIS-HASTINGS ALGORITHM

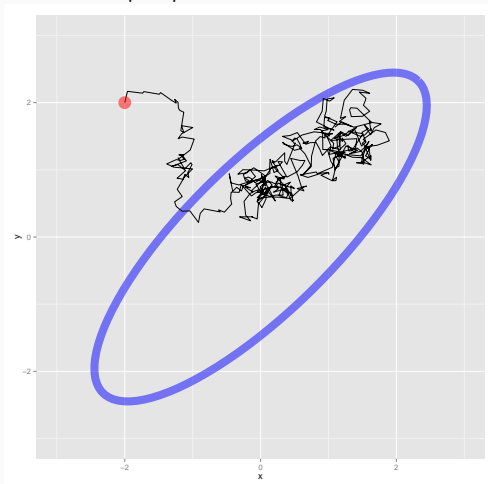
How do we chose the proposal variance?



$$\sigma^2 = 1$$

THE METROPOLIS-HASTINGS ALGORITHM

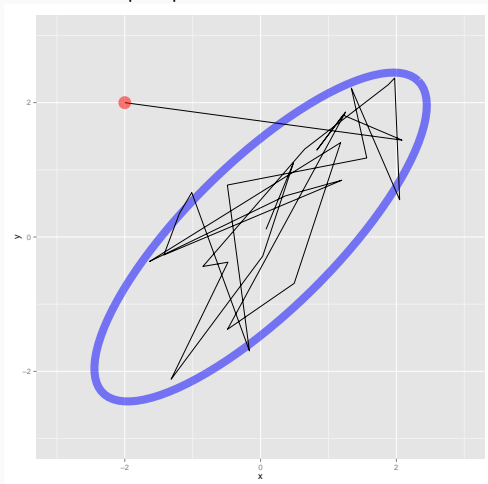
How do we chose the proposal variance?



$$\sigma^2 = .1$$

THE METROPOLIS-HASTINGS ALGORITHM

How do we chose the proposal variance?



$$\sigma^2 = 5$$

DOES THIS SATISFY DETAILED BALANCE?

First, what is the transition kernel $\mathcal{T}(x_{n+1}|x_n)$?

DOES THIS SATISFY DETAILED BALANCE?

First, what is the transition kernel $\mathcal{T}(x_{n+1}|x_n)$?

Prob. of moving from x_n to x_{n+1} is $\alpha(x_{n+1}, x_n)q(x_{n+1}|x_n)$.

DOES THIS SATISFY DETAILED BALANCE?

First, what is the transition kernel $\mathcal{T}(x_{n+1}|x_n)$?

Prob. of moving from x_n to x_{n+1} is $\alpha(x_{n+1}, x_n)q(x_{n+1}|x_n)$.

Prob. of accepting move at x_n is $\alpha(x) = \int_{\mathcal{X}} \alpha(y, x_n)q(y|x_n)dy$.

DOES THIS SATISFY DETAILED BALANCE?

First, what is the transition kernel $\mathcal{T}(x_{n+1}|x_n)$?

Prob. of moving from x_n to x_{n+1} is $\alpha(x_{n+1}, x_n)q(x_{n+1}|x_n)$.

Prob. of accepting move at x_n is $\alpha(x) = \int_{\mathcal{X}} \alpha(y, x_n)q(y|x_n)dy$.

Prob. of rejection at x_n is $r(x_n) = 1 - \alpha(x_n)$.

DOES THIS SATISFY DETAILED BALANCE?

First, what is the transition kernel $\mathcal{T}(x_{n+1}|x_n)$?

Prob. of moving from x_n to x_{n+1} is $\alpha(x_{n+1}, x_n)q(x_{n+1}|x_n)$.

Prob. of accepting move at x_n is $\alpha(x) = \int_{\mathcal{X}} \alpha(y, x_n)q(y|x_n)dy$.

Prob. of rejection at x_n is $r(x_n) = 1 - \alpha(x_n)$.

We then have:

$$\mathcal{T}(x_{n+1}|x_n) = \alpha(x_{n+1}, x_n)q(x_{n+1}|x_n) + r(x_n)\delta(x_n = x_{n+1})$$

DOES THIS SATISFY DETAILED BALANCE?

First, what is the transition kernel $\mathcal{T}(x_{n+1}|x_n)$?

Prob. of moving from x_n to x_{n+1} is $\alpha(x_{n+1}, x_n)q(x_{n+1}|x_n)$.

Prob. of accepting move at x_n is $\alpha(x) = \int_{\mathcal{X}} \alpha(y, x_n)q(y|x_n)dy$.

Prob. of rejection at x_n is $r(x_n) = 1 - \alpha(x_n)$.

We then have:

$$\mathcal{T}(x_{n+1}|x_n) = \alpha(x_{n+1}, x_n)q(x_{n+1}|x_n) + r(x_n)\delta(x_n = x_{n+1})$$

We want to show detailed balance:

$$\pi(x_n)\mathcal{T}(x_{n+1}|x_n) = \pi(x_{n+1})\mathcal{T}(x_n|x_{n+1})$$

THE METROPOLIS-HASTINGS ALGORITHM

Detailed balance: $\pi(x_n)\mathcal{T}(x_{n+1}|x_n) = \pi(x_{n+1})\mathcal{T}(x_n|x_{n+1})$

Consider the LHS:

$$\pi(x_n)\mathcal{T}(x_{n+1}|x_n) = \pi(x_n) (\alpha(x_{n+1}, x_n)q(x_{n+1}, x_n) + r(x_n)\delta(x_n = x_{n+1}))$$

THE METROPOLIS-HASTINGS ALGORITHM

Detailed balance: $\pi(x_n)\mathcal{T}(x_{n+1}|x_n) = \pi(x_{n+1})\mathcal{T}(x_n|x_{n+1})$

Consider the LHS:

$$\pi(x_n)\mathcal{T}(x_{n+1}|x_n) = \pi(x_n) (\alpha(x_{n+1}, x_n)q(x_{n+1}, x_n) + r(x_n)\delta(x_n = x_{n+1}))$$

The first term is:

$$\begin{aligned} & \frac{f(x_n)}{Z} \min \left(1, \frac{f(x_{n+1})q(x_n|x_{n+1})}{f(x_n)q(x_{n+1}|x_n)} \right) q(x_{n+1}|x_n) \\ &= \frac{1}{Z} \min (f(x_n)q(x_{n+1}|x_n), f(x_{n+1})q(x_n|x_{n+1})) \end{aligned}$$

THE METROPOLIS-HASTINGS ALGORITHM

Detailed balance: $\pi(x_n)\mathcal{T}(x_{n+1}|x_n) = \pi(x_{n+1})\mathcal{T}(x_n|x_{n+1})$

Consider the LHS:

$$\pi(x_n)\mathcal{T}(x_{n+1}|x_n) = \pi(x_n) (\alpha(x_{n+1}, x_n)q(x_{n+1}, x_n) + r(x_n)\delta(x_n = x_{n+1}))$$

The first term is:

$$\begin{aligned} & \frac{f(x_n)}{Z} \min \left(1, \frac{f(x_{n+1})q(x_n|x_{n+1})}{f(x_n)q(x_{n+1}|x_n)} \right) q(x_{n+1}|x_n) \\ &= \frac{1}{Z} \min (f(x_n)q(x_{n+1}|x_n), f(x_{n+1})q(x_n|x_{n+1})) \end{aligned}$$

The second term takes this form too.

THE METROPOLIS-HASTINGS ALGORITHM

Detailed balance: $\pi(x_n)\mathcal{T}(x_{n+1}|x_n) = \pi(x_{n+1})\mathcal{T}(x_n|x_{n+1})$

Consider the LHS:

$$\pi(x_n)\mathcal{T}(x_{n+1}|x_n) = \pi(x_n) (\alpha(x_{n+1}, x_n)q(x_{n+1}, x_n) + r(x_n)\delta(x_n = x_{n+1}))$$

The first term is:

$$\begin{aligned} & \frac{f(x_n)}{Z} \min \left(1, \frac{f(x_{n+1})q(x_n|x_{n+1})}{f(x_n)q(x_{n+1}|x_n)} \right) q(x_{n+1}|x_n) \\ &= \frac{1}{Z} \min (f(x_n)q(x_{n+1}|x_n), f(x_{n+1})q(x_n|x_{n+1})) \end{aligned}$$

The second term takes this form too. Thus

$$\pi(x_n)\mathcal{T}(x_{n+1}|x_n) = \pi(x_{n+1})\mathcal{T}(x_n|x_{n+1})$$

Consider a Markov chain on over a set of variables (x_1, \dots, x_d) .

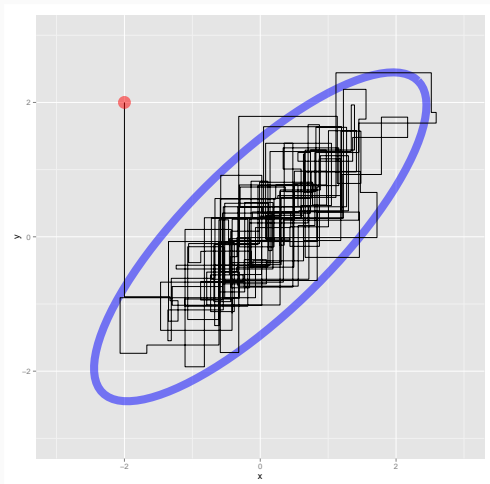
Gibbs sampling cycles through these sequentially (or randomly).

At the i th step, it updates x_i conditioned on the the rest:

$$x_i \sim \pi(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) = \pi(x_i | \mathbf{x}_{\setminus i})$$

Often these conditionals have a much simpler form than the joint.

GIBBS SAMPLING



DETAILED BALANCE FOR GIBBS SAMPLER

Suppose we update component i with prob. ρ_i . Let \mathbf{x} and \mathbf{x}' differ only in component i . Then:

$$\mathcal{T}(\mathbf{x}'|\mathbf{x}) = \rho_i \pi(x'_i | \mathbf{x}_{\setminus i})$$

Also

$$\begin{aligned}\pi(\mathbf{x})\mathcal{T}(\mathbf{x}'|\mathbf{x}) &= \pi(\mathbf{x})\rho_i\pi(x'_i|\mathbf{x}_{\setminus i}) \\ &= \pi(\mathbf{x}_{\setminus i})\pi(x_i|\mathbf{x}_{\setminus i})\rho_i\pi(x'_i|\mathbf{x}_{\setminus i})\end{aligned}$$

From symmetry (or by calculating RHS), we have detailed balance.

Under mild conditions, Gibbs sampling is irreducible.

DETAILED BALANCE FOR GIBBS SAMPLER

Under mild conditions, Gibbs sampling is irreducible.

Can break down under constraints. E.g. two perfectly coupled variables.

Performance deteriorates with strong coupling between variables.

Poor mixing due to coupled variables is always a concern.

DETAILED BALANCE FOR GIBBS SAMPLER

Under mild conditions, Gibbs sampling is irreducible.

Can break down under constraints. E.g. two perfectly coupled variables.

Performance deteriorates with strong coupling between variables.

Poor mixing due to coupled variables is always a concern.

Advantages: Simple, with no free parameters.

DETAILED BALANCE FOR GIBBS SAMPLER

Under mild conditions, Gibbs sampling is irreducible.
Can break down under constraints. E.g. two perfectly coupled variables.
Performance deteriorates with strong coupling between variables.
Poor mixing due to coupled variables is always a concern.

Advantages: Simple, with no free parameters.
Often, conditional independencies in a model along with suitable conjugate priors allow efficient 'blocked-Gibbs samplers'.