

LECTURE 20: OVERVIEW OF MARKOV CHAIN MONTE CARLO

STAT 598Z: INTRODUCTION TO COMPUTING FOR STATISTICS

Vinayak Rao

Department of Statistics, Purdue University

April 16, 2018

Recall Monte Carlo: produce independent samples from $p(x)$, and use sample averages to approximate expectations.

In high dims, hard to even sample from $p(x)$!

Recall Monte Carlo: produce independent samples from $p(x)$, and use sample averages to approximate expectations.

In high dims, hard to even sample from $p(x)$!

Rather than making independent proposals, exploit previous proposals to make good proposals

Allows us to find and explore useful regions of X -space

MARKOV CHAIN MONTE CARLO

Recall Monte Carlo: produce independent samples from $p(x)$, and use sample averages to approximate expectations.

In high dims, hard to even sample from $p(x)$!

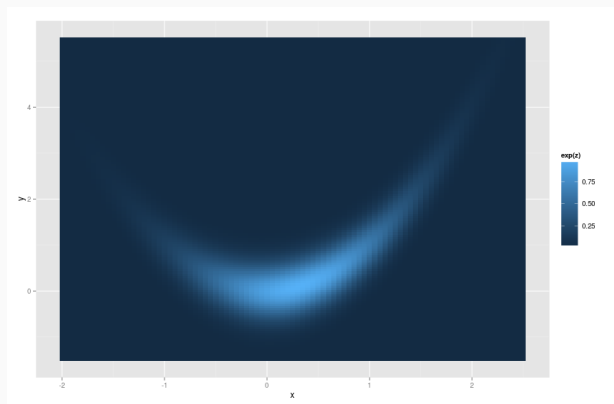
Rather than making independent proposals, exploit previous proposals to make good proposals

Allows us to find and explore useful regions of X -space

Simplest case: use current proposal to make a new proposal

The resulting algorithm: Markov chain Monte Carlo.

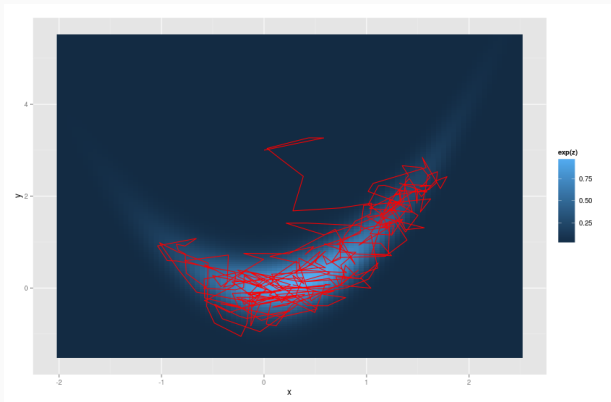
(A Markov chain: future independent of past given present)



The Rosenbrock density (a.k.a. the banana density)

$$p(x, y) \propto \exp \left(-(a - x)^2 - b(y - x^2)^2 \right) \quad (\text{here } a = .3, b = 3)$$

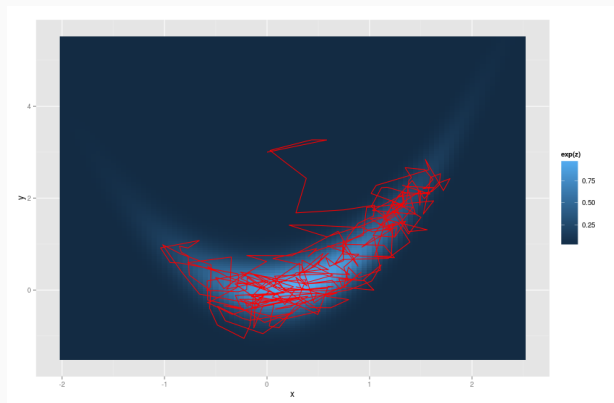
MARKOV CHAIN MONTE CARLO



A random walk:

- start somewhere arbitrary
- make local moves

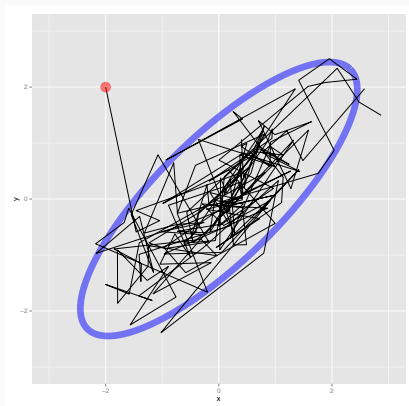
MARKOV CHAIN MONTE CARLO



- Discard initial ‘burn-in’ samples
- Use remaining to obtain Monte Carlo estimates:

$$\frac{1}{N} \sum_{i=1}^N f(x_i) \approx \mathbb{E}_p[f]$$

MARKOV CHAIN MONTE CARLO



A random walk over a 2-d Gaussian

- The goal of MCMC is to find a set of local moves that produce samples (asymptotically) from the right distribution

- The goal of MCMC is to find a set of local moves that produce samples (asymptotically) from the right distribution
- The art of MCMC is to find local moves that converge rapidly (a chain that 'mixes rapidly')

Let element x_i of the chain have distribution p_i

Write $T(x_i \rightarrow x_{i+1}) = p(x_{i+1}|x_i)$ for the transition kernel of the chain.

Then, the $(i + 1)$ st element has distribution

$$p_{i+1}(x_{i+1}) = \int T(x_i \rightarrow x_{i+1})p_i(x_i)dx_i$$

p is the **stationary/equilibrium** distribution of the Markov chain if

$$p(x') = \int T(x \rightarrow x')p(x)dx$$

MCMC: A FIRST LOOK

For a transition function $T(\cdot \rightarrow \cdot)$ with stationary distribution p

- Initialize x_0 from some distribution p_0
- Run a Markov chain for $(B + N)$ iterations with transition T

All x_i for $i > B$ are approximately distributed as p

MCMC: A FIRST LOOK

For a transition function $T(\cdot \rightarrow \cdot)$ with stationary distribution p

- Initialize x_0 from some distribution p_0
- Run a Markov chain for $(B + N)$ iterations with transition T

All x_i for $i > B$ are approximately distributed as p

- Discard the first B ‘burn-in’ samples
- Calculate Monte Carlo average with remaining N samples

$$\frac{1}{N} \sum_{i=B+1}^{B+N} f(x_i) \approx \mathbb{E}_p[f]$$

MCMC: A FIRST LOOK

For a transition function $T(\cdot \rightarrow \cdot)$ with stationary distribution p

- Initialize x_0 from some distribution p_0
- Run a Markov chain for $(B + N)$ iterations with transition T

All x_i for $i > B$ are approximately distributed as p

- Discard the first B 'burn-in' samples
- Calculate Monte Carlo average with remaining N samples

$$\frac{1}{N} \sum_{i=B+1}^{B+N} f(x_i) \approx \mathbb{E}_p[f]$$

Markov chain Monte Carlo estimate

We want to sample from a probability distribution $p(x) = \frac{f(x)}{Z}$

How do we design an appropriate transition kernel $T(\cdot \rightarrow \cdot)$?

We want to sample from a probability distribution $p(x) = \frac{f(x)}{Z}$

How do we design an appropriate transition kernel $T(\cdot \rightarrow \cdot)$?

Different MCMC algorithms take different approaches

The simplest is the Metropolis-Hastings algorithm

Metropolis-Hastings (MH):

- Let current state be x_i
- Propose a new state from $q(w|x_i)$

If we set $x_{i+1} = w$ the resulting Markov chain will have the wrong stationary distribution

Metropolis-Hastings (MH):

- Let current state be x_i
- Propose a new state from $q(w|x_i)$

If we set $x_{i+1} = w$ the resulting Markov chain will have the wrong stationary distribution

- Instead, set $x_{i+1} = w$ (accept) with probability

$$\min\left(1, \frac{p(w)q(x_i|w)}{p(x_i)q(w|x_i)}\right) = \min\left(1, \frac{f(w)q(x_i|w)}{f(x_i)q(w|x_i)}\right)$$

Otherwise, set $x_{i+1} = x_i$ (reject)

Metropolis-Hastings (MH):

- Let current state be x_i
- Propose a new state from $q(w|x_i)$

If we set $x_{i+1} = w$ the resulting Markov chain will have the wrong stationary distribution

- Instead, set $x_{i+1} = w$ (accept) with probability

$$\min\left(1, \frac{p(w)q(x_i|w)}{p(x_i)q(w|x_i)}\right) = \min\left(1, \frac{f(w)q(x_i|w)}{f(x_i)q(w|x_i)}\right)$$

Otherwise, set $x_{i+1} = x_i$ (reject)

Under mild conditions, this corrected Markov chain has the right stationary distribution

METROPOLIS-HASTINGS

Works for any choice of q so long as it's possible to get from any part of space to any other (eventually)

Works for any choice of q so long as it's possible to get from any part of space to any other (eventually)

Acceptance probability:

$$\min\left(1, \frac{p(y)q(x_i|y)}{p(x_i)q(y|x_i)}\right) = \min\left(1, \frac{f(y)q(x_i|y)}{f(x_i)q(y|x_i)}\right)$$

We just have to evaluate the target density $p(x) = \frac{f(x)}{Z}$ up to a proportionality constant

Don't need the normalization constant Z !

Works for any choice of q so long as it's possible to get from any part of space to any other (eventually)

Acceptance probability:

$$\min\left(1, \frac{p(y)q(x_i|y)}{p(x_i)q(y|x_i)}\right) = \min\left(1, \frac{f(y)q(x_i|y)}{f(x_i)q(y|x_i)}\right)$$

We just have to evaluate the target density $p(x) = \frac{f(x)}{Z}$ up to a proportionality constant

Don't need the normalization constant Z !

We only need to

- sample from q
- evaluate transition probabilities $q(y|x)$
- evaluate the target density upto a normalization constant

CHOICE OF PROPOSAL DISTRIBUTION q

- Common choice is a Gaussian centered at previous sample:

$$w|x_i \sim \mathcal{N}(x_i, \sigma^2)$$

- Equivalently,

$$w = x_i + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

CHOICE OF PROPOSAL DISTRIBUTION q

- Common choice is a Gaussian centered at previous sample:

$$w|x_i \sim \mathcal{N}(x_i, \sigma^2)$$

- Equivalently,

$$w = x_i + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

For this proposal $q(w|x_i) = q(x_i|w)$

$$\min\left(1, \frac{f(w)q(x_i|w)}{f(x_i)q(w|x_i)}\right) = \min\left(1, \frac{f(w)}{f(x_i)}\right)$$

CHOICE OF PROPOSAL DISTRIBUTION q

- Common choice is a Gaussian centered at previous sample:

$$w|x_i \sim \mathcal{N}(x_i, \sigma^2)$$

- Equivalently,

$$w = x_i + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

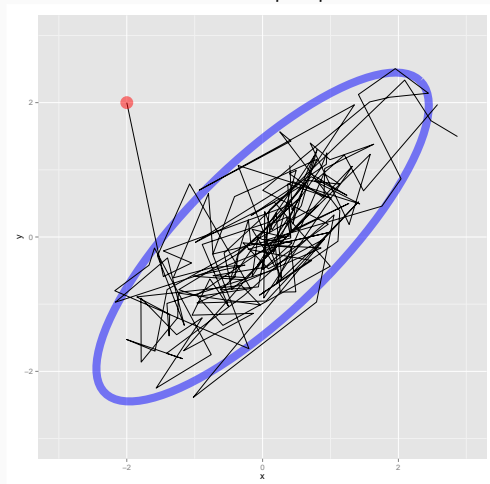
For this proposal $q(w|x_i) = q(x_i|w)$

$$\min\left(1, \frac{f(w)q(x_i|w)}{f(x_i)q(w|x_i)}\right) = \min\left(1, \frac{f(w)}{f(x_i)}\right)$$

- always accept better proposals
- sometimes accept worse proposals

THE METROPOLIS-HASTINGS ALGORITHM

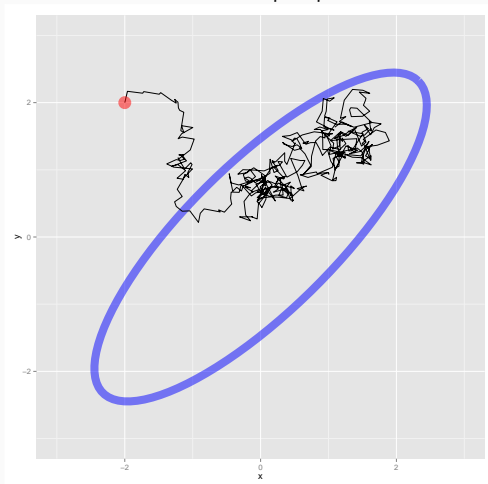
How do we choose the proposal variance?



$$\sigma^2 = 1$$

THE METROPOLIS-HASTINGS ALGORITHM

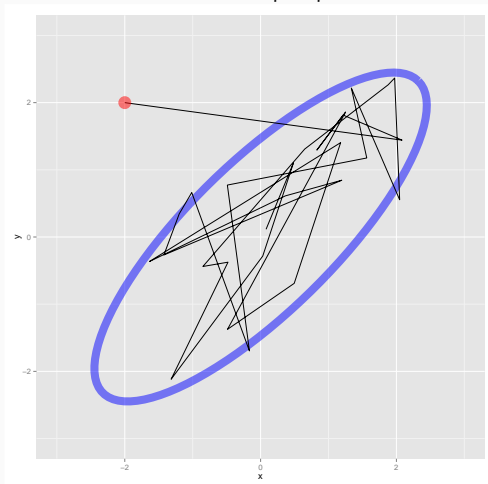
How do we choose the proposal variance?



$$\sigma^2 = .1$$

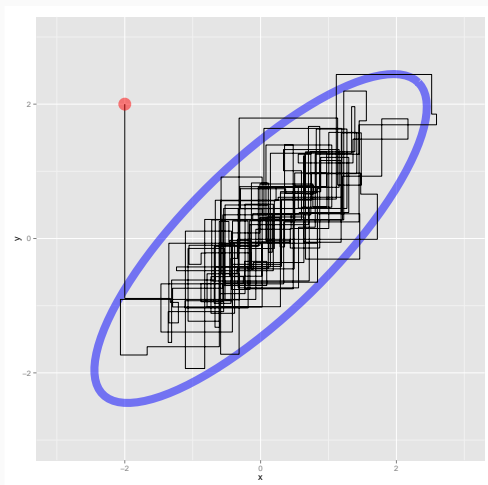
THE METROPOLIS-HASTINGS ALGORITHM

How do we choose the proposal variance?



$$\sigma^2 = 5$$

Sample one component at a time



Consider 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:

$$x(i) \sim p(x(i)|x(1), \dots, x(i-1), x(i+1), \dots, x(n))$$

Often these 1-d conditionals are much simpler than the joint

Think of coordinate descent