

Session 4

Bayesian Methods

Yuri Balasanov

University of Chicago, MScA

© Y. Balasanov, iLykei 2016

© Yuri Balasanov, iLykei 2016

All Rights Reserved

No part of this lecture notes document or any of its contents may be reproduced, copied, modified or adapted without the prior written consent of the author, unless otherwise indicated for stand-alone materials.

The content of these lectures, any comments, opinions and materials are put together by the author especially for the course Linear and Nonlinear Statistical Models, they are sole responsibility of the author, but not of the author's employers or clients.

The author cannot be held responsible for any material damage as a result of use of the materials presented in this document or in this course.

For any inquiries contact the author, Yuri Balasanov, at
ybalasan@uchicago.edu or yuri.balasanov@iLykei.com .

Outline of the Session

- Motivation for Monte Carlo methods in Bayesian analysis
- History of Monte Carlo and MCMC
- Metropolis algorithm
- Gibbs sampling
- Summary of Metropolis and Gibbs sampling methods

Limitations of Bayesian Approach without Monte Carlo Techniques

- We tried Bayesian analysis using analytical approach with conjugate distributions and grid approximation.
Both methods have limitations:
 - Analytical approach is limited by the number of conjugate pairs
 - Grid approximation is limited by dimensionality of the parameter space
- Application of Markov Chain Monte Carlo (MCMC) allows us to extend Bayesian analysis to wide range of important practical applications. This has been a breakthrough of the last couple decades
- The main idea of MCMC is approximation of posterior distribution by sampling from it
- MCMC requires that prior distribution $p(\theta)$ can be easily computed for any θ and that likelihood function can also be efficiently calculated for any sample and any θ

Monte Carlo and MCMC

Who Started it?



Figure: Missing Pictures: Arianna Rosenbluth, Augusta H. Teller.

Motivation for MCMC

- Grid approximation stops working in multidimensional parameter space because it requires N^d points evaluated, where N is the number of grid points for one parameter and d is the number of parameters
- Monte Carlo is expected to reduce the number of sampling points. But if simulation is done independently then it also requires the number of calculations that grows exponentially with growing number of dimensions d
- The only solution that may work is if simulations are not independent. In this case random walk needs to find the most populated area quickly and then spend more time exploring it
- This leads to a random walk in the form of Markov process: next step depends on the current state, but independent from the past

Metropolis Algorithm: Simple Case

Let parameter θ take values $\{\theta_1, \dots, \theta_k\}$ with probabilities $\{p_1 = p(\theta_1), \dots, p_k = p(\theta_k)\}$.

The steps of Metropolis algorithm are:

- 1 Let the current simulated value be θ_c .
- 2 The next candidate for simulation is selected between θ_{c-1} and θ_{c+1} with probabilities 0.5. Let the next candidate be θ_n .
- 3 Then if $p(\theta_n) > p(\theta_c)$ the next simulated value is θ_n .
If $p(\theta_n) \leq p(\theta_c)$ then there is another step of randomization: the next simulated value is θ_n with probability $p = \frac{p(\theta_n)}{p(\theta_c)}$ and θ_c repeats with probability $(1 - p)$.

After repeating such steps N times the probabilities of simulated values $\{\theta_1, \dots, \theta_k\}$ will be very close to $\{p_1 = p(\theta_1), \dots, p_k = p(\theta_k)\}$

Remarkable properties of MCMC?

Note that we needed to do three simple things in order to implement MCMC:

- 1 Generate the next proposed value θ_n
- 2 Calculate $p(\theta_n)$ and $\frac{p(\theta_n)}{p(\theta_c)}$
- 3 Simulate uniform random variable $\zeta \in [0, 1]$

The real value of MCMC is realized when probability distribution $p(\theta | \text{Data})$ is posterior distribution and it does not have to be normalized, meaning that we just can use $p(\theta | \text{Data}) \propto p(\text{Data} | \theta) p(\theta)$

Markov Chain in Metropolis Algorithm

Metropolis algorithm: Markov chain.

$p(i|j) = P\{\theta_i | \theta_j\}$ probability of moving from θ_j to θ_i . Note: $p(i|i) = 0; |i-j| > 1$

$$p(i|i+1) = 0.5 \min\left[\frac{P(i+1)}{P(i)}, 1\right] \quad p(i+1|i) = 0.5 \min\left[\frac{P(i)}{P(i+1)}, 1\right]$$

$$\frac{p(i|i+1)}{p(i+1|i)} = \frac{0.5 \min\left[\frac{P(i+1)}{P(i)}, 1\right]}{0.5 \min\left[\frac{P(i)}{P(i+1)}, 1\right]} = \begin{cases} \frac{1}{\frac{P(i)}{P(i+1)}}, & p(i+1) > p(i) \\ \frac{\frac{P(i+1)}{P(i)}}{1}, & p(i+1) < p(i) \end{cases} = \frac{P(i+1)}{P(i)}$$

Ratio of transition probabilities back and forth equals ratio of target probabilities.

Transition probabilities depend only on previous step: Markov property.

Transition Matrix in Metropolis Algorithm

Transition Matrix.

$$\begin{bmatrix} p(1|1) & p(2|1) & 0 & \dots & \dots & \dots \\ p(1|2) & p(2|2) & p(3|2) & 0 & \dots & \dots \\ 0 & p(2|3) & p(3|3) & p(4|3) & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \dots & 0 & p(k-1|k) & p(k|k) & p(k+1|k) & 0 \dots \\ \dots & 0 & \dots & \dots & p(n-1|n) & p(n|n) & p(1|n) \end{bmatrix}$$

This is a matrix of transition probabilities of a Markov chain H .

Transitions over q steps are given by H^q .

Such Markov chain has stable distribution and stable distribution is the target distribution

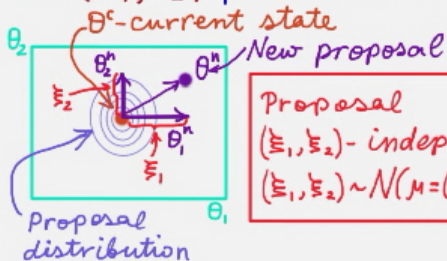
Thus convergence of the process !!!

General Form of Metropolis Algorithm I

- Metropolis algorithm was published in 1953 by Metropolis, Rosenbluth, Rosenbluth Teller & Teller and got the name from the first author. Some sources give credit for the algorithm to his coauthors Marshall and Arianna Rosenbluth
- General algorithm applies to continuous multidimensional parameter space with posterior distribution $p(\theta|D)$, $\theta = \{\theta_1, \dots, \theta_k\}$
- It starts from θ^0 , for which $p(\theta^0|D) \neq 0$, $\theta^0 = \{\theta_1^0, \dots, \theta_k^0\}$
- If current point is θ^c then next candidate θ^{new} is generated from some simple distribution. For example, following slide shows random walk in 2-parametric state space using 2-dimensional normal distribution with mean $(0,0)$, standard deviation $(0.2, 0.2)$ and zero correlation
- Once the next candidate is selected the decision is made the same way as in multinomial Metropolis example: candidate θ^{new} is accepted with probability 1 if $p(\theta^{new}|D) > p(\theta^c|D)$; if $p(\theta^{new}|D) \leq p(\theta^c|D)$ then θ^{new} is accepted with probability $\frac{p(\theta^{new}|D)}{p(\theta^c|D)}$, otherwise the next sampling point is θ_c again

General Form of Metropolis Sampling II

Metropolis sampling
 $\theta = (\theta_1, \theta_2)$ parameter space.



Proposal

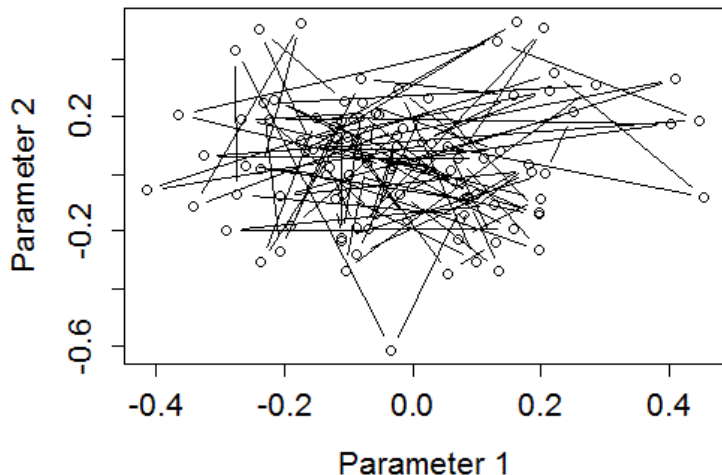
(ξ_1, ξ_2) - independent sample

$$(\xi_1, \xi_2) \sim N(\mu = (0, 0), \Sigma = (\sigma_1, \sigma_2), \rho = 0)$$

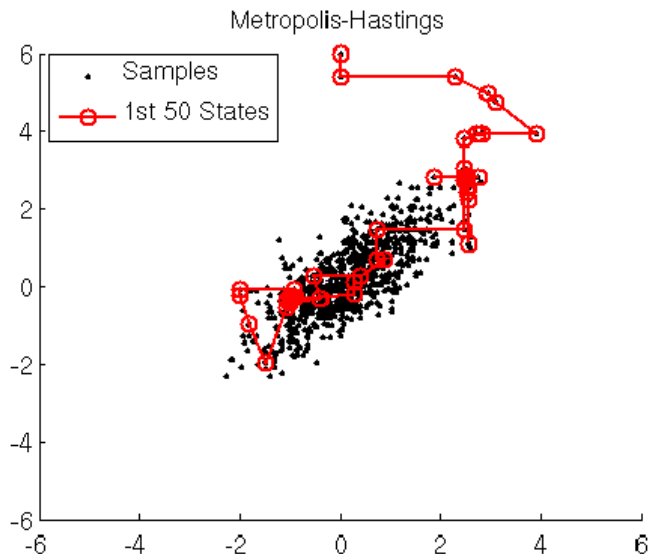
Acceptance: $p(\theta^n | \mathcal{D}) > p(\theta^c | \mathcal{D}) \Rightarrow P\{\theta^n\} = 1$

$$p(\theta^n | \mathcal{D}) \leq p(\theta^c | \mathcal{D}) \Rightarrow \begin{cases} P\{\theta^n\} = \frac{p(\theta^n | \mathcal{D})}{p(\theta^c | \mathcal{D})} \\ P\{\theta^c\} = 1 - \frac{p(\theta^n | \mathcal{D})}{p(\theta^c | \mathcal{D})} \end{cases}$$

Random Walk for 2 Parameters



Metropolis-Hastings Random Walk: Exploration of Distribution



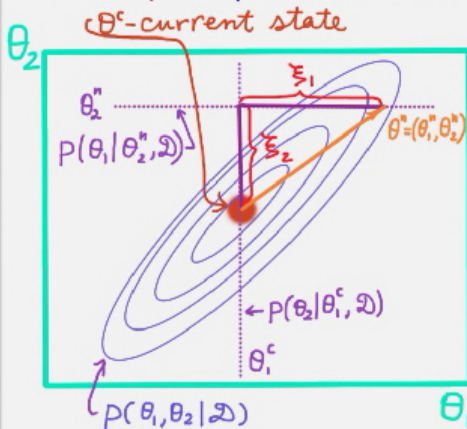
Gibbs Sampling

- Metropolis sampling is a very general method replacing straight Monte Carlo with simulations of a Markov chain process that is efficient numerically and converges to the target distribution faster
- A significant limitation of Metropolis MCMC appears in multidimensional parameters space: next candidate is generated by an arbitrary proposal distribution
- If the proposal distribution is very different from the target distribution then significant number of proposed candidates is rejected reducing efficiency of the algorithm
- A modification of the Metropolis MCMC algorithm samples from one-dimensional posterior marginal distribution directly solving the problem of tuning proposal distribution to the posterior distribution. This modification is called **Gibbs sampling**

Gibbs Sampling II

Gibbs Sampling.

$\theta = (\theta_1, \theta_2)$ parameter space.



Step 1. Generate $\xi_2 \sim p(\theta_2 | \theta_1^c, \mathcal{D})$; $\theta_2^n \leftarrow \xi_2$

Step 2. Generate $\xi_1 \sim p(\theta_1 | \theta_2^n, \mathcal{D})$; $\theta_1^n \leftarrow \xi_1$

Step 3. Define proposal: $\theta^n = (\theta_1^n, \theta_2^n)$

Decide between θ^c and θ^n using the same rule as Metropolis

Summary of the Two Approaches

- Metropolis algorithm is a general algorithm for sampling from a given distribution. Distribution does not have to be normalized.
- In multidimensional state space proposal distribution needs to be tuned to the posterior distribution to avoid multiple rejections
- Gibbs sampling is a particular case of Metropolis algorithm in which generation of a new proposal consists of steps of sampling from marginal distributions for each of the parameters
- Gibbs sampling does not require tuning to posterior distribution, it samples directly from it. But it requires knowing the posterior distribution
- Both distributions are very efficient for analysis of hierarchical models