

Gibbs Sampling

Sargur Srihari

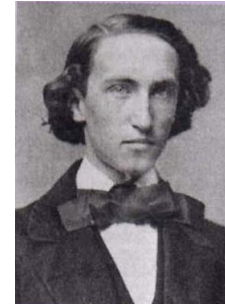
srihari@cedar.buffalo.edu

Topics

1. Gibbs Sampling
2. Use with graphical models
3. Gibbs as a case of Metropolis Hastings
4. Proof of sampling
5. Slice sampling

Gibbs Sampling

- A simple and widely applicable MCMC algorithm
 - Special case of Metropolis-Hastings
- Consider distribution $p(\mathbf{z}) = p(z_1, \dots, z_M)$ from which we wish to sample
- We have chosen an initial state for the Markov chain
- Each step involves replacing value of one variable by a value drawn from $p(z_i | \mathbf{z}_{\setminus i})$
 where symbol $\mathbf{z}_{\setminus i}$ denotes z_1, \dots, z_M with z_i omitted
- Repeat procedure by cycling through variables in some particular order



Josiah Willard Gibbs

1839-1903

Born New Haven CT

First US PhD in Engg.

Developed
vector analysis,
crystallography and
planetary orbits

Gibbs sampling was
named as such by
Geman and Geman
In 1984

When is Gibbs sampling used?

- Gibbs sampling is applicable when the joint distribution is not known explicitly or is difficult to sample from directly, but the conditional distribution of each variable is known and is easy (or at least, easier) to sample from
 - This is usually true in PGMs

Gibbs Sampling Usage

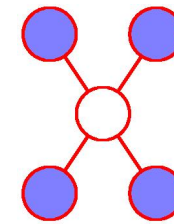
- Gibbs Sampling is an MCMC that samples each random variable of a PGM, one at a time
 - Gibbs is a special case of the MH algorithm
- Gibbs Sampling algorithms...
 - Are fairly easy to derive for many graphical models
 - e.g. mixture models, Latent Dirichlet allocation
 - Have reasonable computation and memory requirements, since they sample one r.v. at a time
 - Can be Rao-Blackwellized (integrate out some r.v.s)
 - to decrease sampling variance – called collapsed Gibbs sampling

The Gibbs Sampling Algorithm

- 1. Suppose the PGM contains variables x_1, \dots, x_n
- 2. Initialize starting values for x_1, \dots, x_n
- 3. Do until convergence:
 - 1. Pick an ordering of the n variables (can be fixed or random)
 - 2. For each variable x_i in order:
 - 1. Sample $x \sim p(x_i \mid x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$,
 - i.e. the conditional distribution of x_i given the current values of all other variables
 - 2. Update $x_i \leftarrow x$
- When we update x_i , we immediately use its new value for sampling other variables x_j

Gibbs sampling and Graphical Models

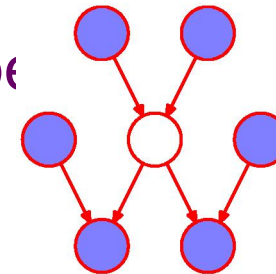
- Practical Applicability of Gibbs sampling



Markov blanket
for undirected graph

- depends on ease with which samples can be drawn from $p(z_k | z_{\setminus k})$
- In the case of PGMs

- Conditional distributions for nodes depend only on variables in Markov blanket
 - which are its neighbors in the graph



Markov blanket
for directed graph

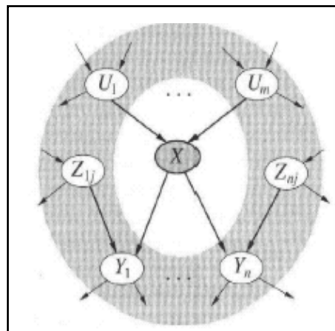
- Gibbs sampling is a distributed algorithm
 - It is not parallel since samples generated sequentially

Markov Blanket to the rescue

- The conditional $p(x_i \mid x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$ looks intimidating, but recall Markov Blankets
- Let $MB(x_i)$ be the Markov Blanket of x_i , then

$$p(x_i \mid x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) = P(x_i \mid MB(x_i))$$

- For a BN, the Markov Blanket of x_i is the set containing its parents, children, and co-parents



- For an MRF, the Markov Blanket of x_i is its immediate neighbors

Gibbs as a case of M-H

- Gibbs sampling can be seen as a special case of MH with proposal

$$q(x'_i, x_{-i} / x_i, x_{-i}) = p(x'_i / x_{-i})$$

- It is easy check that the acceptance probability simplifies to one
- Assuming the right transition operator, both Gibbs sampling and MH will eventually produce samples from their stationary distribution, which by construction is p

Gibbs with Three Variables

- Distribution $p(z_1, z_2, z_3)$ over three variables
- At step t selected values are $z_1^{(t)}$, $z_2^{(t)}$ and $z_3^{(t)}$
- Replace $z_1^{(t)}$ by new value $z_1^{(t+1)}$ obtained by sampling from $p(z_1 | z_2^{(t)}, z_3^{(t)})$
- Replace $z_2^{(t)}$ by value $z_2^{(t+1)}$ by sampling from $p(z_2 | z_1^{(t+1)}, z_3^{(t)})$
 - New value for z_1 is used straightaway
- Update z_3 with a sample $z_3^{(t+1)}$ drawn from $p(z_3 | z_1^{(t+1)}, z_2^{(t+1)})$
- Cycle through three variables in turn

Gibbs Sampling with M variables

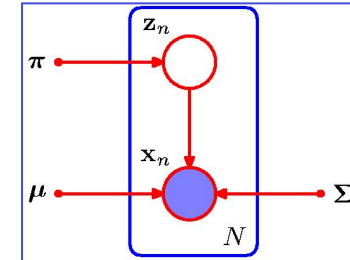
- Initialize first sample: $\{z_i, i=1, \dots, M\}$
- For $t=1, \dots, T$, $T = \text{no of samples}$
 - Sample $z_1^{(t+1)} \sim p(z_1 | z_2^{(t)}, z_3^{(t)}, \dots, z_M^{(t)})$
 - Sample $z_2^{(t+1)} \sim p(z_2 | z_1^{(t+1)}, z_3^{(t)}, \dots, z_M^{(t)})$
 -
 - Sample $z_j^{(t+1)} \sim p(z_j | z_1^{(t+1)}, \dots, z_{j-1}^{(t+1)}, z_{j+1}^{(t)}, \dots, z_M^{(t)})$
 -
 - Sample $z_M^{(t+1)} \sim p(z_M | z_1^{(t+1)}, z_2^{(t+1)}, \dots, z_{M-1}^{(t+1)})$
- $p(z_j | z_{-j})$ is called a *full conditional* for variable j

Gibbs Sampling for inferring parameters of GMM

- Full joint distribution

$$\begin{aligned}
 p(\mathbf{x}, \mathbf{z}, \pi, \mu, \Sigma) &= p(\mathbf{x} | \mathbf{z}, \mu, \pi, \Sigma) p(\mathbf{z}, \pi, \mu, \Sigma) \\
 &= p(\mathbf{x} | \mathbf{z}, \mu, \Sigma) p(\mathbf{z}, \pi) p(\mu, \Sigma) ; \text{ given } \mathbf{z} \text{ we don't need } \pi \\
 &= p(\mathbf{x} | \mathbf{z}, \mu, \Sigma) p(\mathbf{z} | \pi) p(\pi) p(\mu) p(\Sigma)
 \end{aligned}$$

– Using a semi-conjugate prior



$$\begin{aligned}
 p(\mathbf{x}, \mathbf{z}, \pi, \mu, \Sigma) &= p(\mathbf{x} | \mathbf{z}, \mu, \Sigma) p(\mathbf{z} | \pi) p(\pi) \prod_{k=1}^K p(\mu_k) p(\Sigma_k) \\
 &= \left(\prod_{n=1}^N \prod_{k=1}^K \left(\pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right)^{I(z_n=k)} \right) \times Dir(\pi | \alpha) \prod_{k=1}^K N(\mu_k | m_0, V_0) IW(\Sigma_k | S_0, \nu_0)
 \end{aligned}$$

– IW is the inverse Wishart distribution

- We need full conditionals (to obtain samples) for

- Discrete indicators, $p(z_n = l | \mathbf{x}_n, \mu, \pi, \Sigma)$

- Mixing weights $p(\pi | \mathbf{z})$ for π_1, \dots, π_K

- Means $p(\mu_k | \Sigma_k, \mathbf{z}, \mathbf{x})$

- Covariances $p(\Sigma_k | \mu_k, \mathbf{z}, \mathbf{x})$

- Note that we know the values of \mathbf{x}_n

Full Conditionals for Gibbs GMM

- Discrete Indicators

$$p(z_n = k \mid \mathbf{x}_n, \pi, \mu, \Sigma) \propto \pi_k N(\mathbf{x}_n \mid \mu_k, \Sigma_k)$$

- Mixing weights

$$p(\pi \mid \mathbf{z}) = \text{Dir} \left(\left\{ \alpha_k + \sum_{n=1}^N I(z_n = k) \right\}_{k=1}^K \right)$$

- Means

$$p(\mu_k \mid \Sigma_k, \mathbf{z}, \mathbf{x}) = N(\mu_k \mid \mathbf{m}_k, V_k)$$

Terms V_k and S_k are sample statistics

- Covariances

$$p(\Sigma_k \mid \mu_k, \mathbf{z}, \mathbf{x}) = IW(\Sigma_k \mid S_k, \nu_k)$$

Proof of Sampling

- To show that procedure samples from given distribution
- First show that distribution $p(z)$ is invariant (or stationary) of each sampling step and hence of whole Markov chain
- Second requirement is ergodicity
 - Every state reachable from every other state
- The two requirements are formally defined next

Markov chain properties

- Marginal probability for a variable
 - Expressed in terms of marginal probability of previous variable in chain

$$p(z^{(m+1)}) = \sum_{z^{(m)}} p(z^{(m+1)} | z^{(m)}) p(z^{(m)})$$

Probability of sample is sum of probs over all values of prev sample

- Invariant or stationary

- A distribution is invariant wrt a Markov chain if each step leaves the marginal distribution invariant

- $p^*(z)$ is invariant if

$$p^*(z) = \sum_{z'} T(z', z) p^*(z')$$

$$T(z', z) = p(z|z')$$

- Required distribution $p(z)$ is invariant if it satisfies property of detailed balance

$$p^*(z)T(z, z') = p^*(z')T(z', z)$$

Two directions between z and z' are the same

Markov chain that respects detailed balance is reversible.

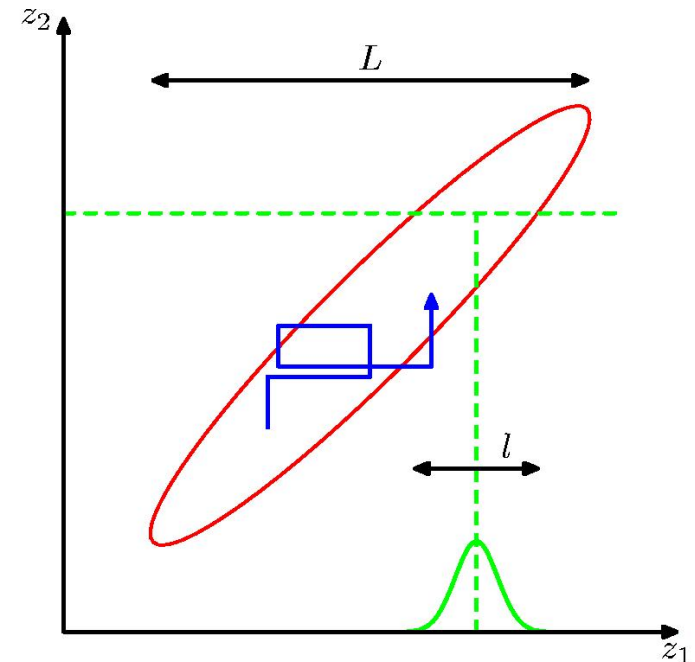
Ergodicity

- Our goal is to use Markov chains to sample from a given distribution
- We need to set up a Markov chain such that the desired distribution is invariant
 - Also, irrespective of choice of initial distribution $p(\mathbf{z}^{(0)})$,
 - As $m \rightarrow \infty$ the distribution $p(\mathbf{z}^{(m)})$ converges to the required invariant distribution $p^*(\mathbf{z})$
- This property is called *ergodicity*
 - And the invariant distribution is called the *equilibrium distribution*

Ergodic also means no state has a zero probability of exit from it
And every state is reachable from every other state

Gibbs sampling of two variables

- Two correlated Gaussian variables
- Step size is governed by standard deviation of conditional distribution
 - Is $O(l)$
 - Leads to slow progress in direction of elongation of the joint distribution
- No of steps needed to obtain an independent sample is $O((L/l)^2)$



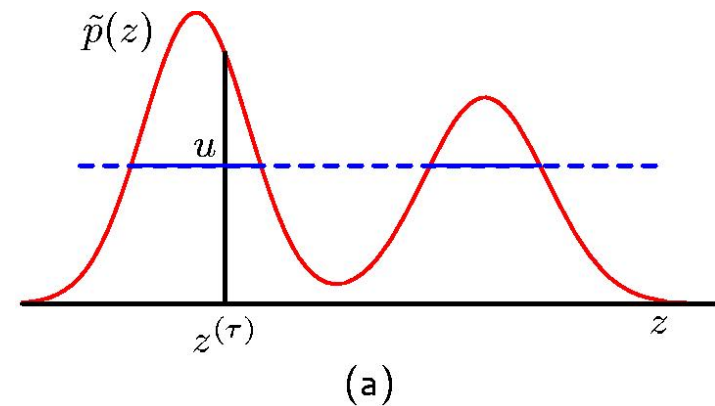
Conditional distribution of width l
Marginal distribution of width L

Slice Sampling

- Metropolis algorithm is sensitive to step size
 - Too small: slow decorrelation due to random walk behavior
 - Too large: inefficiency due to high rejection rate
- Slice sampling provides an adaptive step size to match the distribution

Illustration of Slice

- For given $z^{(t)}$, a value of u is chosen uniformly in region $0 \leq u \leq \tilde{p}(z^{(t)})$
 - which is a slice through the distribution



- Since infeasible to sample from slice, a new sample is drawn from

$$z_{\min} \leq z \leq z_{\max}$$

which contains the previous value $z^{(t)}$

