# Approx. Inference via Markov Chain Monte Carlo

Piyush Rai

Topics in Probabilistic Modeling and Inference (CS698X)

March 27, 2018
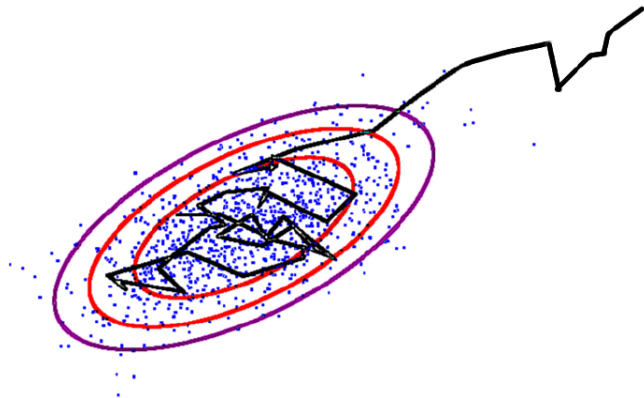
# Markov Chain Monte Carlo (MCMC)

- Goal: Generate samples from some target distribution $p(z) = \frac{\tilde{p}(z)}{Z}$, where $z$ is high-dimensional

- Assume we can evaluate $p(z)$ at least up to a proportionality constant (i.e., can compute $\tilde{p}(z)$)

- Basic idea: MCMC uses a Markov Chain which, when converged, starts giving samples from $p(z)$

$$\underbrace{z^{(1)} \to z^{(2)} \to z^{(3)} \to}_{\text{initial samples typically garbage}} \cdots \to \underbrace{z^{(L-2)} \to z^{(L-1)} \to z^{(L)}}_{\text{after convergence, actual samples from } p(z)}$$

- Given a current sample $z^{(\ell)}$ from the chain, MCMC generates the next sample $z^{(\ell+1)}$ as
  - Use a proposal distribution $q(z|z^{(\ell)})$ to generate a candidate sample $z^*$
  - Accept/reject $z^*$ as the next sample based on an acceptance criterion (will see later)
  - If accepted, $z^{(\ell+1)} = z^*$. If rejected, $z^{(\ell+1)} = z^{(\ell)}$

- Note that in MCMC, the proposal distribution $q(z|z^{(\ell)})$ depends on the previous sample (unlike methods such as rejection sampling)

# Markov Chain Monte Carlo (MCMC)

# MCMC: The Basic Scheme

- MCMC chain run <span style="color:red">infinitely long</span> (i.e., post-convergence) will give ONE sample from the target $p(z)$
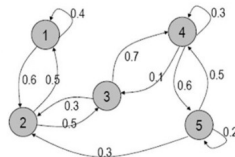


- But we usually require <u>several samples</u> to approximate $p(z)$. How do we get those?
  - Start at an initial $z^{(0)}$. Using a prop. dist. $q(z^{(\ell+1)}|z^{(\ell)})$, run the chain long enough, say $T_1$ steps
  - Discard the first $(T_1 - 1)$ samples (called <span style="color:red">"burn-in" samples</span>) and take the last sample $z^{(T_1)}$
  - Continue from $z^{(T_1)}$ up to $T_2$ steps, discard intermediate samples, take the last sample $z^{(T_2)}$
    - This helps ensure that $z^{(T_1)}$ and $z^{(T_2)}$ are <span style="color:blue">uncorrelated</span>
  - Repeat the same for a total of $S$ times
  - In the end, we have $S$ i.i.d. samples from $p(z)$, i.e., <span style="color:blue">$z^{(T_1)}, z^{(T_2)}, \ldots, z^{(T_S)} \sim p(z)$</span>
  - Note: Good choices for $T_1$ and $T_i - T_{i-1}$ are usually based on heuristics
  - Note: MCMC is an <span style="color:red">approximate method</span> because we don't usually know what $T_1$ is "long enough"

## MCMC: Some Basic Theory

- A first order Markov Chain assumes $p(z^{(\ell+1)}|z^{(1)}, \ldots, z^{(\ell)}) = p(z^{(\ell+1)}|z^{(\ell)})$

- A 1st order Markov Chain $z^{(0)}, z^{(1)}, \ldots, z^{(L)}$ is a sequence of r.v.'s and is defined by the following
  - An initial state distribution $p(z^{(0)})$
  - A Transition Function (TF): $T_\ell(z^{(\ell)} \to z^{(\ell+1)}) = p(z^{(\ell+1)}|z^{(\ell)})$.

- TF defines a distribution over the values of next state given the value of the current state

- Assuming a discrete state-space, the TF is defined by a $K \times K$ probability table

Transition probabilities can be defined using a $K \times K$ table if **z** is a discrete r.v. with $K$ possible values

$$\begin{array}{c} \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{ccccc} 1 & 2 & 3 & 4 & 5 \\ \begin{bmatrix} 0.4 & 0.6 & 0.0 & 0.0 & 0.0 \\ 0.5 & 0.0 & 0.5 & 0.0 & 0.0 \\ 0.0 & 0.3 & 0.0 & 0.7 & 0.0 \\ 0.0 & 0.0 & 0.1 & 0.3 & 0.6 \\ 0.0 & 0.3 & 0.0 & 0.5 & 0.2 \end{bmatrix} \end{array}$$
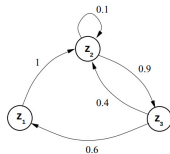


- Homogeneous Markov Chain: The TF is the same for all $\ell$, i.e., $T_\ell = T$

## MCMC: Some Basic Theory

- Consider the following simple TF

$$T = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0.1 & 0.9 \\ 0.6 & 0.4 & 0 \end{bmatrix}$$



- Consider the initial state distribution $p(\mathbf{z}^{(0)}) = [0.5, 0.2, 0.3]$
- Easy to see that $p(\mathbf{z}^{(0)}) \times T = [0.2, 0.6, 0.2] \Rightarrow$ distribution of $\mathbf{z}^{(1)}$
- Also easy to see that, after a few (say $m$) iterations, $p(\mathbf{z}^{(0)}) \times T^m = [0.2, 0.4, 0.4] = p(\mathbf{z})$ (say)
- For the above $T$, <u>any choice</u> of $p(\mathbf{z}^{(0)})$ leads to $p(\mathbf{z})$
  - Such a $p(\mathbf{z})$ is called the stationary/invariant distribution of this Markov Chain
- A Markov Chain has a stationary distribution if $T$ has the following properties
  - Irreducibility: $T$'s graph is connected (ensures reachability from anywhere to anywhere)
  - Aperiodicity: $T$'s graph has no cycles (ensures that the chain isn't trapped in cycles)

## MCMC: Some Basic Theory

- A Markov Chain has a stationary distribution $p(z)$ is the chain satisfies detailed balance

$$p(z)T(z \rightarrow z') = p(z')T(z' \rightarrow z)$$

- Integrating out (or summing over) both sides w.r.t. $z'$ gives

$$p(z) = \int p(z')T(z' \rightarrow z)dz'$$

- Therefore $p(z)$ is a stationary distribution of this chain

- Thus a Markov Chain with detailed balance will always converge to a stationary distribution

# Some MCMC Algorithms

# Metropolis-Hastings (MH) Sampling (Hastings, 1970)

- Suppose we wish to generate samples from a distribution $p(z) = \frac{\tilde{p}(z)}{Z_p}$

- Assume a proposal distribution $q(z|z^{(\tau)})$, e.g., $\mathcal{N}(z|z^{(\tau)}, \sigma^2 I_D)$

- In each step, draw $z^* \sim q(z|z^{(\tau)})$ and accept the sample $z^*$ with probability

$$A(z^*, z^{(\tau)}) = \min\left(1, \frac{\tilde{p}(z^*)q(z^{(\tau)}|z^*)}{\tilde{p}(z^{(\tau)})q(z^*|z^{(\tau)})}\right)$$

- The acceptance probability makes intuitive sense. Note the kind of $z^*$ would it favor/unfavor:
  - It favors accepting $z^*$ if $\tilde{p}(z^*)$ has a higher value than $\tilde{p}(z^{(\tau)})$
  - Unfavors $z^*$ if the proposal distribution $q$ unduly favors it (i.e., if $q(z^*|z^{(\tau)})$ is large)
  - Favors $z^*$ if we can "reverse" to $z^{(\tau)}$ from $z^*$ (i.e., if $q(z^{(\tau)}|z^*)$ is large). Needed for good "mixing"

- Transition function of this Markov Chain: $T(z^{(\tau)} \rightarrow z^*) = A(z^*, z^{(\tau)})q(z^*|z^{(\tau)})$

- Exercise: Show that $T(z \rightarrow z^{(\tau)})$ satisfies the detailed balance property

$$T(z \rightarrow z^{(\tau)})p(z) = T(z^{(\tau)} \rightarrow z)p(z^{(\tau)})$$

## The MH Sampling Algorithm

- Initialize $z^{(0)}$ randomly

- For $\ell = 0, \ldots, L - 1$

  - Sample $u \sim \text{Unif}(0, 1)$

  - Sample $z^* \sim q(z^* | z^{(\ell)})$

  - If $u < A(z^*, z^{(\ell)}) = \min \left( 1, \frac{\tilde{p}(z^*) q(z^{(\ell)} | z^*)}{\tilde{p}(z^{(\ell)}) q(z^* | z^{(\ell)})} \right)$

    $$z^{(\ell+1)} = z^* \qquad \text{(meaning: accepting with probability } A(z^*, z^{(\ell)}))$$
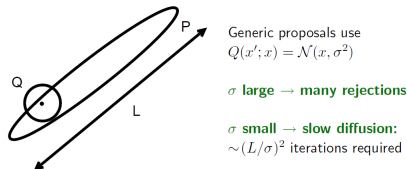
    else

    $$z^{(\ell+1)} = z^{(\ell)}$$

# MH Sampling: Some Comments

- If proposal distrib. is symmetric, we get Metropolis Sampling algorithm (Metropolis, 1953) with

$$A(\boldsymbol{z}^*, \boldsymbol{z}^{(\tau)}) = \min\left(1, \frac{\tilde{p}(\boldsymbol{z}^*)}{\tilde{p}(\boldsymbol{z}^{(\tau)})}\right)$$

- Some limitations of MH sampling

  - MH can have a very slow convergence. Figure below: $P$ is the target dist., $Q$ is the proposal



Generic proposals use
$Q(x'; x) = \mathcal{N}(x, \sigma^2)$

$\sigma$ large → many rejections

$\sigma$ small → slow diffusion:
$\sim (L/\sigma)^2$ iterations required

  - Computing acceptance probability can be expensive: When $p(\boldsymbol{z}) = \frac{\tilde{p}(\boldsymbol{z})}{Z_p}$ represents a posterior distribution of some model, $\tilde{p}$ is the unnormalized posterior that depends on all the data (note: a lot of recent work on speeding up this step using subsets of data)

# Gibbs Sampling (Geman & Geman, 1984)

- Suppose we wish to sample from a joint distribution $p(\mathbf{z})$ where $\mathbf{z} = (z_1, z_2, \ldots, z_M)$
- However, suppose we can't sample from $p(\mathbf{z})$ but can sample from each conditional $p(z_i|\mathbf{z}_{-i})$
  - Can we done easily if we have a locally conjugate model
- Gibbs sampling uses the conditionals $p(z_i|\mathbf{z}_{-i})$ as the proposal distribution
- Gibbs sampling samples from these conditionals in a cyclic order
- Gibbs sampling is equivalent to Metropolis Hastings sampling with acceptance prob. $= 1$

$$A(\mathbf{z}^*, \mathbf{z}) = \frac{p(\mathbf{z}^*)q(\mathbf{z}|\mathbf{z}^*)}{p(\mathbf{z})q(\mathbf{z}^*|\mathbf{z})} = \frac{p(z_i^*|\mathbf{z}_{-i}^*)p(\mathbf{z}_{-i}^*)p(z_i|\mathbf{z}_{-i}^*)}{p(z_i|\mathbf{z}_{-i})p(\mathbf{z}_{-i})p(z_i^*|\mathbf{z}_{-i})} = 1$$

where we use the fact that $\mathbf{z}_{-i}^* = \mathbf{z}_{-i}$

## Gibbs Sampling: Sketch of the Algorithm

$M$: Total number of variables, $T$: number of Gibbs sampling steps
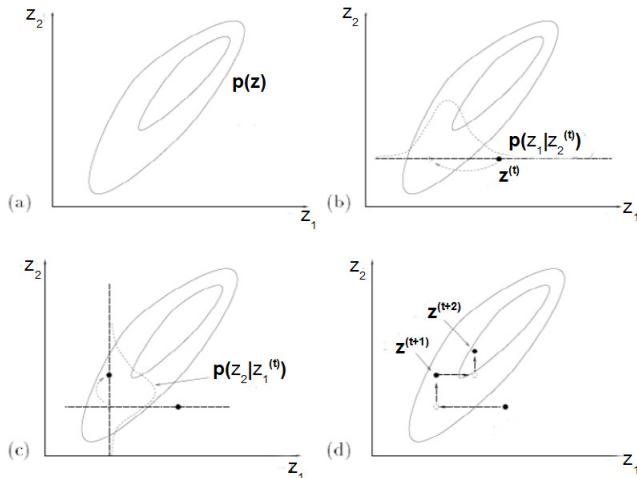
1. Initialize $\{z_i : i = 1, \ldots, M\}$
2. For $\tau = 1, \ldots, T$:
   - Sample $z_1^{(\tau+1)} \sim p(z_1 | z_2^{(\tau)}, z_3^{(\tau)}, \ldots, z_M^{(\tau)})$.
   - Sample $z_2^{(\tau+1)} \sim p(z_2 | z_1^{(\tau+1)}, z_3^{(\tau)}, \ldots, z_M^{(\tau)})$.
     ⋮
   - Sample $z_j^{(\tau+1)} \sim p(z_j | z_1^{(\tau+1)}, \ldots, z_{j-1}^{(\tau+1)}, z_{j+1}^{(\tau)}, \ldots, z_M^{(\tau)})$.
     ⋮
   - Sample $z_M^{(\tau+1)} \sim p(z_M | z_1^{(\tau+1)}, z_2^{(\tau+1)}, \ldots, z_{M-1}^{(\tau+1)})$.

Note: When sampling each variable from its conditional posterior, we use the most recent values of all other variables (this is akin to a co-ordinate ascent like procedure)

Note: Order of updating the variables *usually* doesn't matter (but see "Scan Order in Gibbs Sampling: Models in Which it Matters and Bounds on How Much" from NIPS 2016)

# Gibbs Sampling: A Simple Example

Can sample from a 2-D Gaussian using 1-D Gaussians (recall that if the joint distribution is a 2-D Gaussian, conditionals will simply be 1-D Gaussians)

# Gibbs Sampling: Some Comments

- One of the most popular MCMC algorithm

- Very easy to derive and implement for locally conjugate models

- Many variations exist, e.g.,

  - Blocked Gibbs: sample multiple variables jointly (sometimes possible)
  - Rao-Blackwellized Gibbs: Can collapse (i.e., integrate out) the unneeded variables while sampling. Also called "collapsed" Gibbs sampling
  - MH within Gibbs

- Instead of sampling from the conditionals, an alternative is to use the mode of the conditional.

  - Called the "Itearative Conditional Mode" (ICM) algorithm (doesn't give the posterior though)

# Using MCMC Samples: An Example

- Consider making <u>predictions</u> in the PMF model: $\mathbf{R} = \mathbf{U}\mathbf{V}^\top + \epsilon$, with $\mathbf{R}$ being $N \times M$, $\mathbf{U}$ being $N \times K$, $\mathbf{V}$ being $M \times K$, and each entry $r_{ij} = \boldsymbol{u}_i^\top \boldsymbol{v}_j + \epsilon_{ij}$

- Suppose our MCMC sampler (e.g., Gibbs sampler) gave us $S$ samples $\{\mathbf{U}^{(s)}, \mathbf{V}^{(s)}\}_{s=1}^{S}$

- How to make predictions for a missing $r_{ij}$ (its predictive mean, its predictive variance)?

- Given the samples $\{\mathbf{U}^{(s)}, \mathbf{V}^{(s)}\}_{s=1}^{S}$ from the posterior, we can approximate the mean of $r_{ij}$ as
$$\mathbb{E}[r_{ij}|\mathbf{R}] = \mathbb{E}[\boldsymbol{u}_i^\top \boldsymbol{v}_j + \epsilon_{ij}] = \mathbb{E}[\boldsymbol{u}_i^\top \boldsymbol{v}_j] \approx \frac{1}{S} \sum_{s=1}^{S} \boldsymbol{u}_i^{(s)\top} \boldsymbol{v}_j^{(s)}$$

- The variance can be likewise approximated as
$$
\begin{aligned}
\mathrm{Var}[r_{ij}|\mathbf{R}] = \mathrm{Var}[\boldsymbol{u}_i^\top \boldsymbol{v}_j + \epsilon_{ij}] &= \mathrm{Var}[\boldsymbol{u}_i^\top \boldsymbol{v}_j] + \mathrm{Var}[\epsilon_{ij}] \\
&= \mathbb{E}[(\boldsymbol{u}_i^\top \boldsymbol{v}_j)^2] - \left[\mathbb{E}[\boldsymbol{u}_i^\top \boldsymbol{v}_j]\right]^2 + \beta^{-1} \\
&\approx \frac{1}{S} \sum_{s=1}^{S} \left(\boldsymbol{u}_i^{(s)\top} \boldsymbol{v}_j^{(s)}\right)^2 - \left(\frac{1}{S} \sum_{s=1}^{S} \boldsymbol{u}_i^{(s)\top} \boldsymbol{v}_j^{(s)}\right)^2 + \beta^{-1}
\end{aligned}
$$

- Question: Can't we average all samples to get a single $\mathbf{U}$ and a single $\mathbf{V}$ and use those?
  - No. Reason: Possible mode or label switching

# Mode/Label Switching

- The posterior of most latent variable models has multiple modes
  - .. even if it is unimodal when the latent variables are <u>known</u>
- Each sampling iteration can give samples of latent variables from one of the modes
- Example: Consider a clustering model like GMM. Likelihood is invariant to label permutations
  - What one sample considers as cluster 1 may become cluster 2 for next sample
- Therefore averaging latent variables or parameters across samples can be meaningless
- Quantities not affected by permutations of latent variables can be safely averaged
  - Probability that two points belong to the same cluster (e.g., in GMM)
  - Predicting the mean/variance of a missing entry $r_{ij}$ in matrix factorization
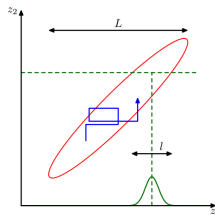
## MCMC/Sampling vs Variational Inference

- MCMC gives samples from the true posterior when run infinitely long

  - But in practice we can never run the chain infinitely long

- VB only gives a locally optimal approximation of the posterior

  - But VB guarantees that we will give it in finite time!

- MCMC is quite general; doesn't assume simplifying assumptions such as mean-field

- Many VB methods (e.g.. BBVI) use sampling as a routine

- MCMC approximation is expensive storage-wise (need to store all the samples)

  - Also, computing any quantity that depends on the posterior is expensive with MCMC because we need to average using all the stored samples!

# MCMC and Random Walk

- MCMC methods use a proposal distribution to draw the next sample given the previous sample

$$\theta^{(t)} \sim \mathcal{N}(\theta^{(t-1)}, \sigma^2)$$

- .. and then we accept/reject (if doing MH) or always accept (if doing Gibbs sampling)

- Such proposal distributions typically lead to a random-walk behavior (e.g., a zig-zag trajectory in Gibbs sampling) and may lead to very slow convergence (pic below: $\theta = [z_1, z_2]$)



- Can be especially critical when the components of $\theta$ are highly correlated

- Using the posterior's gradient may be beneficial for faster convergence (next class)