

Gibbs Sampling Examples, Some Aspects of MCMC

CS698X: Topics in Probabilistic Modeling and Inference

Piyush Rai

Recap: Gibbs Sampling

- An instance of MH sampling where the acceptance probability = 1
- Based on sampling \mathbf{z} one “component” at a time with proposal = conditional distr.

Gibbs Sampling

Initialize $\mathbf{z}^{(0)} = [z_1^{(0)}, z_2^{(0)}, \dots, z_M^{(0)}]$ randomly

For $\ell = 1, \dots, L$

- Sample $\mathbf{z}^{(\ell)}$ by sampling one component at a time (usually cyclic manner)

$\mathbf{z}^{(\ell)}$

$$\begin{aligned} z_1^{(\ell)} &\sim p(z_1 | z_2^{(\ell-1)}, z_3^{(\ell-1)}, \dots, z_M^{(\ell-1)}) \\ z_2^{(\ell)} &\sim p(z_2 | z_1^{(\ell)}, z_3^{(\ell-1)}, \dots, z_M^{(\ell-1)}) \\ &\vdots \\ z_{M-1}^{(\ell)} &\sim p(z_{M-1} | z_1^{(\ell)}, \dots, z_{M-2}^{(\ell)}, z_M^{(\ell-1)}) \\ z_M^{(\ell)} &\sim p(z_M | z_1^{(\ell)}, z_2^{(\ell)}, \dots, z_{M-1}^{(\ell)}) \end{aligned}$$

In practice, we won't use all the L samples to approximate the target distribution $p(\mathbf{z})$ since there will be a burn-in phase and thinning as well



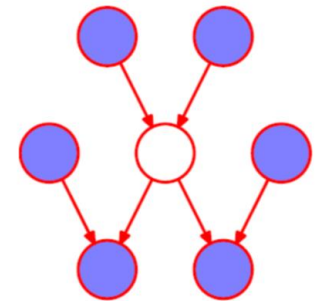
Denoting the collected samples by $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(S)}$, the posterior approximation will be the empirical distribution defined by these samples

- Very easy to derive if the conditional distributions are easy to obtain



Deriving A Gibbs Sampler: The General Recipe

- Suppose the target is an intractable posterior $p(\mathbf{Z}|\mathbf{X})$ where $\mathbf{Z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_M]$
- Gibbs sampling requires the conditional posteriors $p(\mathbf{z}_m|\mathbf{Z}_{-m}, \mathbf{X})$
- In general, $p(\mathbf{z}_m|\mathbf{Z}_{-m}, \mathbf{X}) \propto p(\mathbf{z}_m)p(\mathbf{X}|\mathbf{z}_m, \mathbf{Z}_{-m})$ where \mathbf{Z}_{-m} is assumed “known”
- If $p(\mathbf{z}_m)$ and $p(\mathbf{X}|\mathbf{z}_m, \mathbf{Z}_{-m})$ are conjugate, the above CP is straightforward to obtain
- Another way to get each CP $p(\mathbf{z}_m|\mathbf{Z}_{-m}, \mathbf{X})$ is by following this
 - Write down the expression of $p(\mathbf{X}, \mathbf{Z})$
 - Only terms that contain \mathbf{z}_m needed to get CP of \mathbf{z}_m (up to a prop const)
- In $p(\mathbf{z}_m|\mathbf{Z}_{-m}, \mathbf{X})$, we only need to condition on terms in Markov Blanket of \mathbf{z}_m
 - Markov Blanket of a variable: Its parents, children, and other parents of its children
 - Very useful in deriving CP

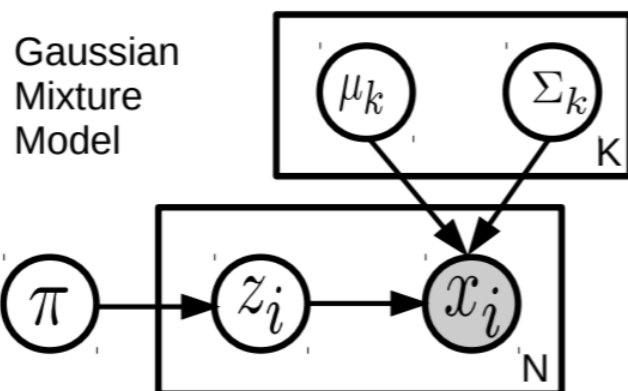


Markov Blanket



Gibbs Sampling: An Example

- The CPs for the Gibbs sampler for a GMM are as shown in green rectangles below



Can verify that Markov Blanket property holds for each CP

Joint distribution of data and unknowns

$$\begin{aligned}
 p(\mathbf{x}, \mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) &= p(\mathbf{x}|\mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Sigma})p(\mathbf{z}|\boldsymbol{\pi})p(\boldsymbol{\pi}) \prod_{k=1}^K p(\boldsymbol{\mu}_k)p(\boldsymbol{\Sigma}_k) \\
 &= \left(\prod_{i=1}^N \prod_{k=1}^K (\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))^{\mathbb{I}(z_i=k)} \right) \times \\
 &\quad \text{Dir}(\boldsymbol{\pi} | \boldsymbol{\alpha}) \prod_{k=1}^K \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_0, \mathbf{V}_0) \text{IW}(\boldsymbol{\Sigma}_k | \mathbf{S}_0, \nu_0)
 \end{aligned}$$

$$p(z_i = k | \mathbf{x}_i, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) \propto \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

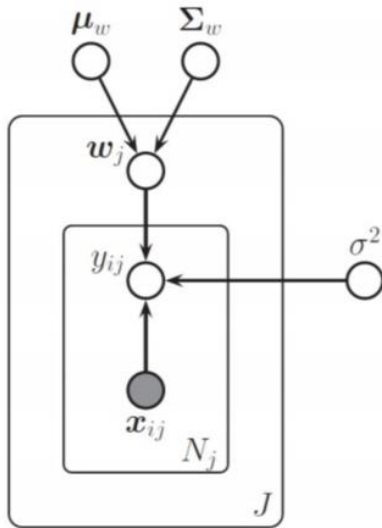
$$p(\boldsymbol{\pi} | \mathbf{z}) = \text{Dir}(\{\alpha_k + \sum_{i=1}^N \mathbb{I}(z_i = k)\}_{k=1}^K)$$

$$\begin{aligned}
 p(\boldsymbol{\mu}_k | \boldsymbol{\Sigma}_k, \mathbf{z}, \mathbf{x}) &= \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_k, \mathbf{V}_k) \\
 \mathbf{V}_k^{-1} &= \mathbf{V}_0^{-1} + N_k \boldsymbol{\Sigma}_k^{-1} \\
 \mathbf{m}_k &= \mathbf{V}_k (\boldsymbol{\Sigma}_k^{-1} N_k \bar{\mathbf{x}}_k + \mathbf{V}_0^{-1} \mathbf{m}_0) \\
 N_k &\triangleq \sum_{i=1}^N \mathbb{I}(z_i = k) \\
 \bar{\mathbf{x}}_k &\triangleq \frac{\sum_{i=1}^N \mathbb{I}(z_i = k) \mathbf{x}_i}{N_k}
 \end{aligned}$$

$$\begin{aligned}
 p(\boldsymbol{\Sigma}_k | \boldsymbol{\mu}_k, \mathbf{z}, \mathbf{x}) &= \text{IW}(\boldsymbol{\Sigma}_k | \mathbf{S}_k, \nu_k) \\
 \mathbf{S}_k &= \mathbf{S}_0 + \sum_{i=1}^N \mathbb{I}(z_i = k) (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \\
 \nu_k &= \nu_0 + N_k
 \end{aligned}$$

Gibbs Sampling: Another Example

J schools
Regression
Problem



$$\begin{aligned}
 & p\left(Y, \{\mathbf{w}_j\}_{j=1}^J, \boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w, \sigma^2 \mid \mathbf{X}\right) \quad \text{Joint distribution of data and unknowns} \\
 &= \left(\prod_{j=1}^J \prod_{i=1}^{N_j} p(y_{ij} \mid \mathbf{x}_{ij}, \mathbf{w}_j, \sigma^2) p(\mathbf{w}_j \mid \boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w) \right) p(\boldsymbol{\mu}_w) p(\boldsymbol{\Sigma}_w) p(\sigma^2) \\
 &= \left(\prod_{j=1}^J \prod_{i=1}^{N_j} \mathcal{N}(y_{ij} \mid \mathbf{w}_j^T \mathbf{x}_{ij}, \sigma^2) \mathcal{N}(\mathbf{w}_j \mid \boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w) \right) \\
 &\quad \mathcal{N}(\boldsymbol{\mu}_w \mid \boldsymbol{\mu}_0, \mathbf{V}_0) \text{IW}(\boldsymbol{\Sigma}_w \mid \boldsymbol{\eta}_0, \mathbf{S}_0^{-1}) \text{IG}(\sigma^2 \mid \nu_0/2, \nu_0 \sigma_0^2/2)
 \end{aligned}$$

Can verify that
Markov Blanket
property holds
for each CP

$$\begin{aligned}
 p(\mathbf{w}_j \mid \mathcal{D}_j, \boldsymbol{\theta}) &= \mathcal{N}(\mathbf{w}_j \mid \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \\
 \boldsymbol{\Sigma}_j^{-1} &= \boldsymbol{\Sigma}^{-1} + \mathbf{X}_j^T \mathbf{X}_j / \sigma^2 \\
 \boldsymbol{\mu}_j &= \boldsymbol{\Sigma}_j (\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \mathbf{X}_j^T \mathbf{y}_j / \sigma^2)
 \end{aligned}$$

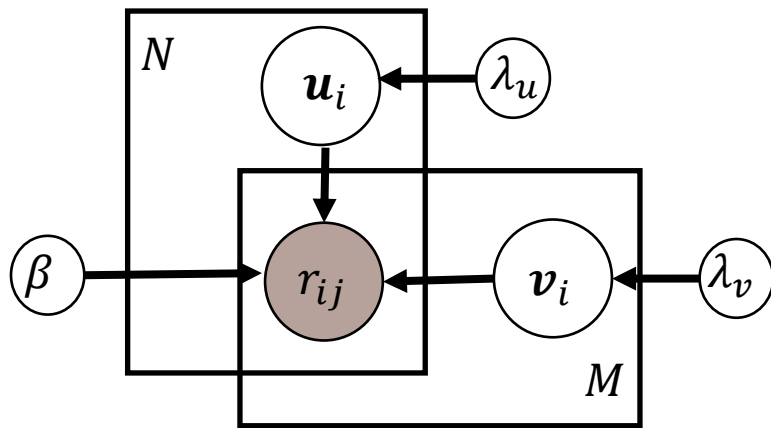
$$\begin{aligned}
 p(\boldsymbol{\mu}_w \mid \mathbf{w}_{1:J}, \boldsymbol{\Sigma}_w) &= \mathcal{N}(\boldsymbol{\mu} \mid \boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N) \\
 \boldsymbol{\Sigma}_N^{-1} &= \mathbf{V}_0^{-1} + J \boldsymbol{\Sigma}^{-1} \\
 \boldsymbol{\mu}_N &= \boldsymbol{\Sigma}_N (\mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + J \boldsymbol{\Sigma}^{-1} \bar{\mathbf{w}}) \\
 \bar{\mathbf{w}} &= \frac{1}{J} \sum_j \mathbf{w}_j
 \end{aligned}$$

$$\begin{aligned}
 p(\boldsymbol{\Sigma}_w \mid \boldsymbol{\mu}_w, \mathbf{w}_{1:J}) &= \text{IW}((\mathbf{S}_0 + \mathbf{S}_\mu)^{-1}, \eta_0 + J) \\
 \mathbf{S}_\mu &= \sum_j (\mathbf{w}_j - \boldsymbol{\mu}_w)(\mathbf{w}_j - \boldsymbol{\mu}_w)^T
 \end{aligned}$$

$$\begin{aligned}
 p(\sigma^2 \mid \mathcal{D}, \mathbf{w}_{1:J}) &= \text{IG}([\nu_0 + N]/2, [\nu_0 \sigma_0^2 + \text{SSR}(\mathbf{w}_{1:J})]/2) \\
 \text{SSR}(\mathbf{w}_{1:J}) &= \sum_{j=1}^J \sum_{i=1}^{N_j} (y_{ij} - \mathbf{w}_j^T \mathbf{x}_{ij})^2
 \end{aligned}$$



Gibbs Sampling: Another Example



Bayesian Matrix Factorization

$$p(\mathbf{R}, \{\mathbf{u}_i\}_{i=1}^N, \{\mathbf{v}_j\}_{j=1}^M, \lambda_u, \lambda_v, \beta)$$

Joint distribution of data and unknowns

Assuming even the hyperparams to be unknown

$$= \prod_{(i,j) \in \Omega} p(r_{ij} | \mathbf{u}_i, \mathbf{v}_j, \beta) \prod_i p(\mathbf{u}_i | \lambda_u) \prod_j p(\mathbf{v}_j | \lambda_v) p(\lambda_u) p(\lambda_v) p(\beta)$$

Can also use non-zero mean and full cov matrix for $\mathbf{u}_i, \mathbf{v}_j$, with Gaussian and Wishart priors respectively*

$$= \prod_{(i,j) \in \Omega} \mathcal{N}(r_{ij} | \mathbf{u}_i^\top \mathbf{v}_j, \beta) \prod_i \mathcal{N}(\mathbf{u}_i | 0, \lambda_u^{-1} \mathbf{I}) \prod_j \mathcal{N}(\mathbf{v}_j | 0, \lambda_v^{-1} \mathbf{I})$$

$$\text{Gamma}(\lambda_u | a, b) \text{Gamma}(\lambda_v | c, d) \text{Gamma}(\beta | e, f)$$

$$p(\mathbf{u}_i | \mathbf{R}, \mathbf{V}) = \mathcal{N}(\mathbf{u}_i | \boldsymbol{\mu}_{u_i}, \boldsymbol{\Sigma}_{u_i})$$

$$\boldsymbol{\Sigma}_{u_i} = (\lambda_u \mathbf{I} + \beta \sum_{j:(i,j) \in \Omega} \mathbf{v}_j \mathbf{v}_j^\top)^{-1}$$

$$\boldsymbol{\mu}_{u_i} = \boldsymbol{\Sigma}_{u_i} (\beta \sum_{j:(i,j) \in \Omega} r_{ij} \mathbf{v}_j)$$

$$p(\mathbf{v}_j | \mathbf{R}, \mathbf{U}) = \mathcal{N}(\mathbf{v}_j | \boldsymbol{\mu}_{v_j}, \boldsymbol{\Sigma}_{v_j})$$

$$\boldsymbol{\Sigma}_{v_j} = (\lambda_v \mathbf{I} + \beta \sum_{i:(i,j) \in \Omega} \mathbf{u}_i \mathbf{u}_i^\top)^{-1}$$

$$\boldsymbol{\mu}_{v_j} = \boldsymbol{\Sigma}_{v_j} (\beta \sum_{i:(i,j) \in \Omega} r_{ij} \mathbf{u}_i)$$

Can verify that Markov Blanket property holds for each CP

$$p(\lambda_u | \mathbf{U}) = \text{Gamma}(\lambda_u | a + 0.5 * NK, b + 0.5 * \sum_{i=1}^N \mathbf{u}_i^\top \mathbf{u}_i)$$

$$p(\lambda_v | \mathbf{V}) = \text{Gamma}(\lambda_v | c + 0.5 * MK, d + 0.5 * \sum_{j=1}^M \mathbf{v}_j^\top \mathbf{v}_j)$$

$$p(\beta | \mathbf{R}, \mathbf{U}, \mathbf{V}) = \text{Gamma}(\beta | e + 0.5 * |\Omega|, f + 0.5 * \sum_{i,j \in \Omega} (r_{ij} - \mathbf{u}_i^\top \mathbf{v}_j)^2)$$

Ω denotes the indices that are observed in the ratings matrix

MCMC: Some Other Aspects



Using the Samples to make Predictions

- Using the S samples $\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(S)}$, our approx. $p(\mathbf{Z}) \approx \frac{1}{S} \sum_{s=1}^S \delta_{\mathbf{Z}^{(s)}}(\mathbf{Z})$

- Any expectation that depends on $p(\mathbf{Z})$ be approximated as

$$\mathbb{E}[f(\mathbf{Z})] = \int f(\mathbf{Z})p(\mathbf{Z})d\mathbf{Z} \approx \frac{1}{S} \sum_{s=1}^S f(\mathbf{Z}^{(s)})$$

- For Bayesian lin. reg., assuming $\mathbf{w}, \beta, \lambda$ to be unknown, the PPD approx. will be

$$\int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) d\mathbf{w} d\beta d\lambda \approx \frac{1}{S} \sum_{s=1}^S p(y_* | \mathbf{x}_*, \mathbf{w}^{(s)}, \beta^{(s)})$$

Joint posterior over all unknowns

Thus, in this case, the PPD is a sum of S Gaussians

Sampling based approximation of PPD

Mean and variance of y_* can be computed using sum of Gaussian properties

Mean: $\mathbb{E}[y_*] = \frac{1}{S} \sum_{s=1}^S \mathbf{w}^{(s)\top} \mathbf{x}_*$

Variance: Exercise! Use definition of variance and use Monte-Carlo approximation

- Sampling based approx. for PPD of other models can also be obtained likewise

Sampling Methods: Label Switching Issue

- Suppose we are given samples $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(S)}$ from the posterior $p(\mathbf{Z}|\mathbf{X})$
- We can't always simply "average" them to get the "posterior mean" $\bar{\mathbf{z}}$
- Why: **Non-identifiability** of latent vars in models with multiple equival. posterior modes
- Example: In clustering via GMM, the likelihood is invariant to how we label clusters
 - What we call cluster 1 in one sample may be cluster 2 in the next sample
 - Say, in GMM, $\mathbf{z}_n^{(1)} = [1, 0]$ and $\mathbf{z}_n^{(2)} = [0, 1]$, both samples imply the same
 - Averaging will give $\bar{\mathbf{z}}_n = [0.5, 0.5]$, which is incorrect
- Quantities not affected by permutations of dims of \mathbf{Z} can be safely averaged
 - E.g., probability that two points belong to the same cluster (e.g., in GMM)
 - Predicting the mean of an entry r_{ij} in matrix factorization $\frac{1}{S} \sum_{s=1}^S \mathbf{u}_i^{(s)\top} \mathbf{v}_j^{(s)}$

One sample may be from near one of the modes and the other may be from near the other mode

Changes in order of entries in these $K \times 1$ vectors across different samples doesn't affect the inner product

MCMC: Some Practical Aspects

- Choice of proposal distribution is important
 - For MH sampling, Gaussian proposal is popular when \mathbf{z} is continuous, e.g.,

$$q(\mathbf{z}|\mathbf{z}^{(\ell-1)}) = \mathcal{N}(\mathbf{z}|\mathbf{z}^{(\ell-1)}, \mathbf{H})$$

Hessian at the MAP of the target distribution

Change at each iter

- Other options: Mixture of proposal distributions, data-driven or adaptive proposals
- Autocorrelation.** Can show that when approximating $f^* = \mathbb{E}[f]$ using $\{\mathbf{Z}^{(s)}\}_{s=1}^S$

$$\bar{f} = \frac{1}{S} \sum_{s=1}^S f_s$$

Monte Carlo assumes uncorrelated samples

Value of f using s^{th} MCMC sample

Basically measures what fractions of the total samples are uncorrelated. Want it to be close to 1

$$\text{var}_{MCMC}[\bar{f}] = \text{var}_{MC}[\bar{f}] + \frac{1}{S^2} \sum_{s \neq t} \mathbb{E}[(f_s - f^*)(f_t - f^*)]$$

$$\text{Effective Sample Size (ESS)} = \frac{\text{var}_{MC}[f]}{\text{var}_{MCMC}[f]}$$

- Autocorrelation function (ACF) at lag t : $\rho_t = \frac{\frac{1}{S-t} \sum_{s=1}^{S-t} (f_s - \bar{f})(f_{s+t} - \bar{f})}{\frac{1}{S-1} \sum_{s=1}^S (f_s - \bar{f})^2}$

Lower is better

- Multiple Chains:** Run multiple chains, take union of generated samples



Approximate Inference: VI vs Sampling

- VI approximates a posterior distribution $p(\mathbf{Z}|\mathbf{X})$ by another distribution $q(\mathbf{Z}|\phi)$
- Sampling uses S samples $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(S)}$ to approximate $p(\mathbf{Z}|\mathbf{X})$
- Sampling can be used within VI (ELBO approx using Monte-Carlo)
- In terms of “comparison” between VI and sampling, a few things to be noted
 - **Convergence:** VI only has local convergence, sampling (in theory) can give exact posterior
 - **Storage:** Sampling based approx needs to storage all samples, VI only needs var. params ϕ
 - **Prediction Cost:** Sampling always requires Monte-Carlo avging for posterior predictive; with VI, sometimes we can get closed form posterior predictive

PPD if using sampling:

$$p(x_*|X) = \int p(x_*|Z)p(Z|X)dZ \approx \frac{1}{S} \sum_{s=1}^S p(x_*|Z^{(s)})$$

PPD if using VI:

$$p(x_*|X) = \int p(x_*|Z)p(Z|X)dZ \approx \int p(x_*|Z)q(Z|\phi)dZ$$

Compressing the S samples into something more compact

- There is some work on “compressing” sampling-based approximations*

*"Compact approximations to Bayesian predictive distributions" by Snelson and Ghahramani, 2005; and "Bayesian Dark Knowledge" by Korattikara et al, 2015



Coming Up Next

- Avoiding the random-walk behavior of MCMC
 - Using gradient information of the posterior
- Scalable MCMC methods

