

Markov Chain Monte Carlo (MCMC) Sampling

CS772A: Probabilistic Machine Learning

Piyush Rai

Markov Chain Monte Carlo (MCMC)

If the target is a posterior, it will be conditioned on data, i.e., $p(\mathbf{z}|\mathbf{x})$

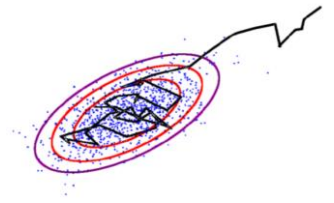
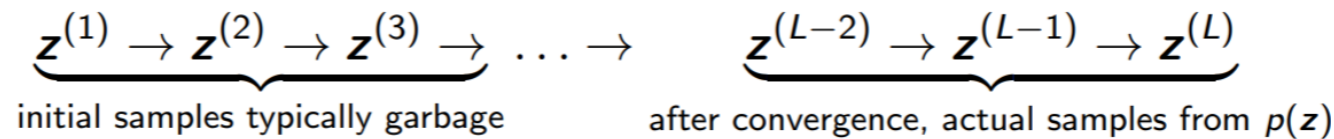
- Goal: Generate samples from some target distribution $p(\mathbf{z}) = \frac{\tilde{p}(\mathbf{z})}{Z_p}$

\mathbf{z} usually is high-dim

- Assume we can evaluate $p(\mathbf{z})$ at least up to a proportionality constant

Means we can at least evaluate $\tilde{p}(\mathbf{z})$

- MCMC uses a **Markov Chain** which, when converged, starts giving samples from $p(\mathbf{z})$



- Given current sample $\mathbf{z}^{(\ell)}$ from the chain, MCMC generates the next sample $\mathbf{z}^{(\ell+1)}$ as

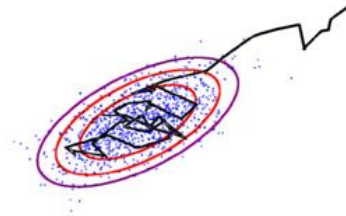
- Use a **proposal distribution** $q(\mathbf{z}|\mathbf{z}^{(\ell)})$ to generate a candidate sample \mathbf{z}_*
- **Accept/reject** \mathbf{z}_* as the next sample based on an **acceptance criterion** (will see later)
- If accepted, set $\mathbf{z}^{(\ell+1)} = \mathbf{z}_*$. If rejected, set $\mathbf{z}^{(\ell+1)}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

Should also have the same support as $p(\mathbf{z})$

- Important: The proposal distribution $q(\mathbf{z}|\mathbf{z}^{(\ell)})$ depends on the previous sample $\mathbf{z}^{(\ell)}$



MCMC: The Basic Scheme



- The chain run infinitely long (i.e., upon convergence) will give ONE sample from $p(\mathbf{z})$
- But we usually require **several samples** to approximate $p(\mathbf{z})$
- This is done as follows
 - Start the chain at an initial $\mathbf{z}^{(0)}$
 - Using the proposal $q(\mathbf{z}|\mathbf{z}^{(\ell)})$, run the chain long enough, say T_1 steps
 - Discard the first $T_1 - 1$ samples (called “**burn-in**” **samples**) and take last sample $\mathbf{z}^{(T_1)}$
 - Continue from $\mathbf{z}^{(T_1)}$ up to T_2 steps, discard intermediate samples, take last sample $\mathbf{z}^{(T_2)}$
 - This discarding (called “**thinning**”) helps ensure that $\mathbf{z}^{(T_1)}$ and $\mathbf{z}^{(T_2)}$ are **uncorrelated**
 - Repeat the same for a total of S times
 - In the end, we now have S *approximately independent* samples from $p(\mathbf{z})$
- Note: Good choices for T_1 and $T_i - T_{i-1}$ (thinning gap) are usually based on heuristics

MCMC is exact in theory but approximate in practice since we can't run the chain for infinitely long in practice

Thus we say that the samples are approximately from the target distribution

Will treat it as our first sample from $p(\mathbf{z})$

Requirement for Monte Carlo approximation

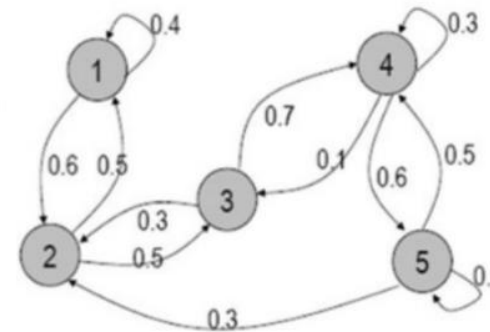


MCMC: Some Basic Theory

- A first order Markov Chain assumes $p(\mathbf{z}^{(\ell+1)} | \mathbf{z}^{(1)}, \dots, \mathbf{z}^{(\ell)}) = p(\mathbf{z}^{(\ell+1)} | \mathbf{z}^{(\ell)})$
- A 1st order Markov Chain $\mathbf{z}^{(0)}, \mathbf{z}^{(1)}, \dots, \mathbf{z}^{(L)}$ is a sequence of r.v.'s and is defined by
 - An initial state distribution $p(\mathbf{z}^{(0)})$
 - A Transition Function (TF): $T_\ell(\mathbf{z}^{(\ell)} \rightarrow \mathbf{z}^{(\ell+1)}) = p(\mathbf{z}^{(\ell+1)} | \mathbf{z}^{(\ell)})$
- TF is a distribution over the values of next state given the value of the current state
- Assuming a K -dim discrete state-space, TF will be $K \times K$ probability table

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

	1	2	3	4	5
1	0.4	0.6	0.0	0.0	0.0
2	0.5	0.0	0.5	0.0	0.0
3	0.0	0.3	0.0	0.7	0.0
4	0.0	0.0	0.1	0.3	0.6
5	0.0	0.3	0.0	0.5	0.2



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



MCMC: Some Basic Theory

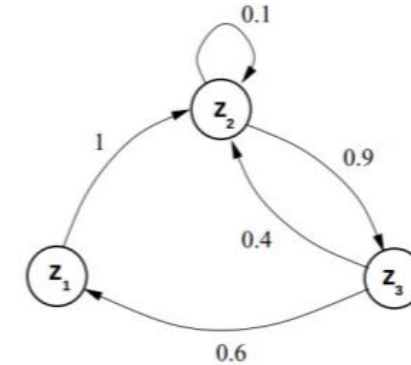
- Consider the following Markov Chain to sample a discrete r.v. \mathbf{z} with 3 possible values

The initial state distribution for \mathbf{z}

$$p(\mathbf{z}^{(0)}) = p(z_1^{(0)}, z_2^{(0)}, z_3^{(0)}) = [0.5, 0.2, 0.3]$$

Probabilities of the initial state taking each of the 3 possible values

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



Distribution of \mathbf{z} after taking the first step

$$p(\mathbf{z}^{(1)}) = p(\mathbf{z}^{(0)}) \times T = [0.2, 0.6, 0.2] \quad (\text{rounded to single digit after decimal})$$

After doing it a few more (say some m) times

$$p(\mathbf{z}^{(0)}) \times T^m = [0.2, 0.4, 0.4] \quad (\text{rounded to single digit after decimal})$$

Stationary/Invariant Distribution $p(\mathbf{z})$ of this Markov Chain

$p(\mathbf{z})$ is multinoulli with $\pi = [0.2, 0.4, 0.4]$

- $p(\mathbf{z})$ being Stationary means no matter what $p(\mathbf{z}^{(0)})$ is, we will reach $p(\mathbf{z})$
- 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 with transition function T has stationary distribution $p(\mathbf{z})$ if T satisfies

Known as the Detailed Balance condition

$$p(\mathbf{z})T(\mathbf{z}'|\mathbf{z}) = p(\mathbf{z}')T(\mathbf{z}|\mathbf{z}')$$

Here $T(b|a)$ denotes the transition probability of going from state a to state b

- Integrating out (or summing over) detailed balanced condition on both sides w.r.t. \mathbf{z}'

Thus $p(\mathbf{z})$ is the stationary distribution of this Markov Chain

$$p(\mathbf{z}) = \int p(\mathbf{z}')T(\mathbf{z}|\mathbf{z}')d\mathbf{z}'$$

- Thus a Markov Chain with detailed balance always converges to a stationary distribution
- Detailed Balance ensures reversibility
- Detailed balance is sufficient but not necessary condition for having a stationary distr.



Some MCMC Algorithms



Metropolis-Hastings (MH) Sampling (1960)

- Suppose we wish to generate samples from a target distribution $p(\mathbf{z}) = \frac{\tilde{p}(\mathbf{z})}{Z_p}$
- Assume a suitable proposal distribution $q(\mathbf{z}|\mathbf{z}^{(\tau)})$, e.g., $\mathcal{N}(\mathbf{z}|\mathbf{z}^{(\tau)}, \sigma^2 \mathbf{I})$
- In each step, draw \mathbf{z}^* from $q(\mathbf{z}|\mathbf{z}^{(\tau)})$ and accept \mathbf{z}^* with probability

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

Favors acceptance of \mathbf{z}^* if it is more probable than $\mathbf{z}^{(\tau)}$ (under $p(\mathbf{z})$)

Downweight the probability of acceptance of \mathbf{z}^* if the proposal itself favors its generation (i.e., if $q(\mathbf{z}^*|\mathbf{z}^{(\tau)})$ is high)

- Transition function of this Markov Chain

- $T(\mathbf{z}^*|\mathbf{z}^{(\tau)}) = A(\mathbf{z}^*, \mathbf{z}^{(\tau)}) q(\mathbf{z}^*|\mathbf{z}^{(\tau)})$ if state changed
- $T(\mathbf{z}^*|\mathbf{z}^{(\tau)}) = q(\mathbf{z}^{(\tau)}|\mathbf{z}^{(\tau)}) + \sum_{\mathbf{z}^* \neq \mathbf{z}^{(\tau)}} (1 - A(\mathbf{z}^*, \mathbf{z}^{(\tau)})) q(\mathbf{z}^*|\mathbf{z}^{(\tau)})$ otherwise

Can Show that this TF satisfied detailed balance



The MH Sampling Algorithm

- Initialize $\mathbf{z}^{(1)}$ randomly
- For $\ell = 1, 2, \dots, L$
 - Sample $\mathbf{z}^* \sim q(\mathbf{z}^* | \mathbf{z}^{(\ell)})$ and $u \sim \text{Unif}(0,1)$
 - Compute acceptance probability

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

- If $A(\mathbf{z}^*, \mathbf{z}^{(\ell)}) > u$

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

Meaning accepting \mathbf{z}^* with probability $A(\mathbf{z}^*, \mathbf{z}^{(\ell)})$

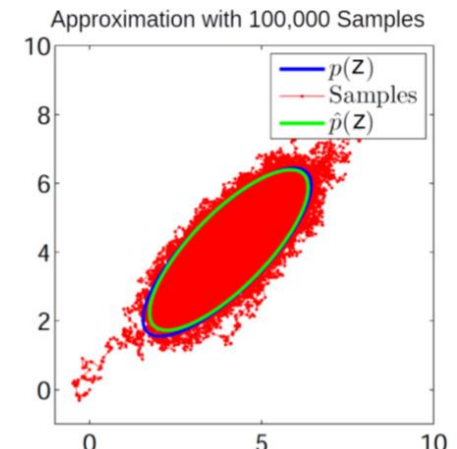
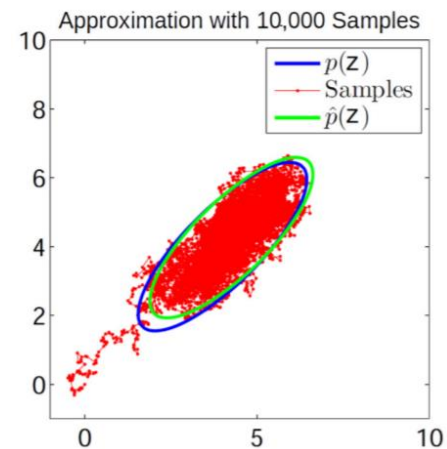
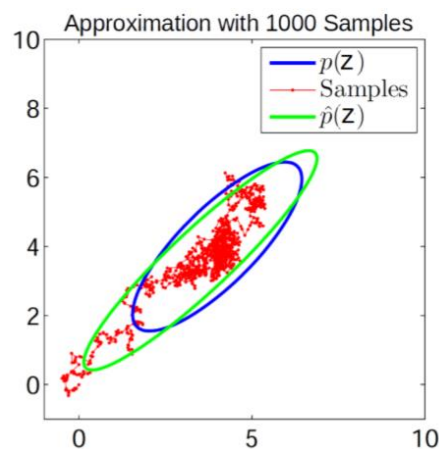
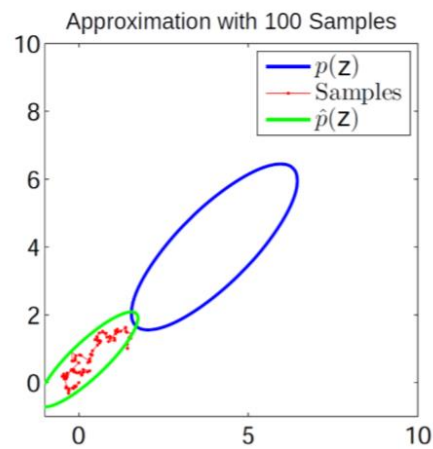
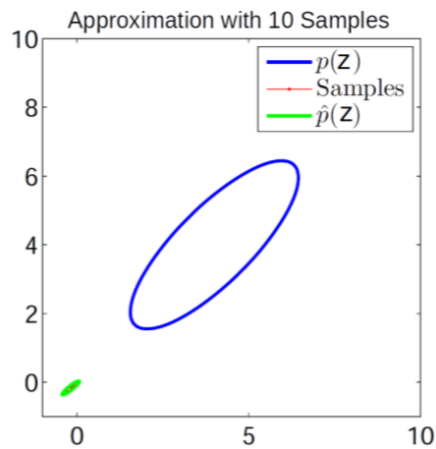
- Else

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



MH Sampling in Action: A Toy Example..

- Target distribution $p(\mathbf{z}) = \mathcal{N} \left(\begin{bmatrix} 4 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix} \right)$
- Proposal distribution $q(\mathbf{z}^{(t)} | \mathbf{z}^{(t-1)}) = \mathcal{N} \left(\mathbf{z}^{(t-1)}, \begin{bmatrix} 0.01 & 0 \\ 0 & 0.01 \end{bmatrix} \right)$

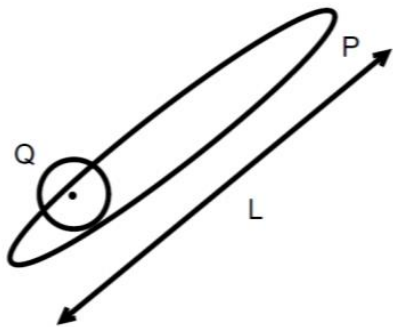


MH Sampling: Some Comments

- If prop. distrib. is symmetric, we get [Metropolis Sampling](#) algo (Metropolis, 1953) with

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

- Some limitations of MH sampling
 - Can sometimes have very slow convergence (also known as slow “mixing”)



$$Q(\mathbf{z}|\mathbf{z}^{(\tau)}) = \mathcal{N}(\mathbf{z}|\mathbf{z}^{(\tau)}, \sigma^2 \mathbf{I})$$

σ large \Rightarrow many rejections

σ small \Rightarrow slow diffusion

$\sim \left(\frac{L}{\sigma}\right)^2$ iterations required for convergence

- Computing acceptance probability can be expensive*, e.g., if $p(\mathbf{z}) = \frac{\tilde{p}(\mathbf{z})}{Z_p}$ is some target posterior then $\tilde{p}(\mathbf{z})$ would require computing likelihood on all the data points (expensive)



Gibbs Sampling (Geman & Geman, 1984)

- Goal: Sample from a joint distribution $p(\mathbf{z})$ where $\mathbf{z} = [z_1, z_2, \dots, z_M]$
- Suppose we can't sample from $p(\mathbf{z})$ but can sample from each conditional $p(z_i | \mathbf{z}_{-i})$
 - In Bayesian models, can be done easily if we have a locally conjugate model
- For Gibbs sampling, the proposal is the conditional distribution $p(z_i | \mathbf{z}_{-i})$
- Gibbs sampling samples from these conditionals in a cyclic order
- Gibbs sampling is equivalent to MH sampling with acceptance prob. = 1

Hence no need to compute it

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

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

Since only one component is changed at a time



Gibbs Sampling: Sketch of the Algorithm

- M : Total number of variables, T : number of Gibbs sampling iterations

1. Initialize $\{z_i : i = 1, \dots, M\}$ Assuming $\mathbf{z} = [z_1, z_2, \dots, z_M]$

2. For $\tau = 1, \dots, T$:

– Sample $z_1^{(\tau+1)} \sim p(z_1 | z_2^{(\tau)}, z_3^{(\tau)}, \dots, z_M^{(\tau)})$.

– Sample $z_2^{(\tau+1)} \sim p(z_2 | z_1^{(\tau+1)}, z_3^{(\tau)}, \dots, z_M^{(\tau)})$.

\vdots

– Sample $z_j^{(\tau+1)} \sim p(z_j | z_1^{(\tau+1)}, \dots, z_{j-1}^{(\tau+1)}, z_{j+1}^{(\tau)}, \dots, z_M^{(\tau)})$.

\vdots

– Sample $z_M^{(\tau+1)} \sim p(z_M | z_1^{(\tau+1)}, z_2^{(\tau+1)}, \dots, z_{M-1}^{(\tau+1)})$.

CP of each component of \mathbf{z} uses the most recent values (from this or the previous iteration) of all the other components

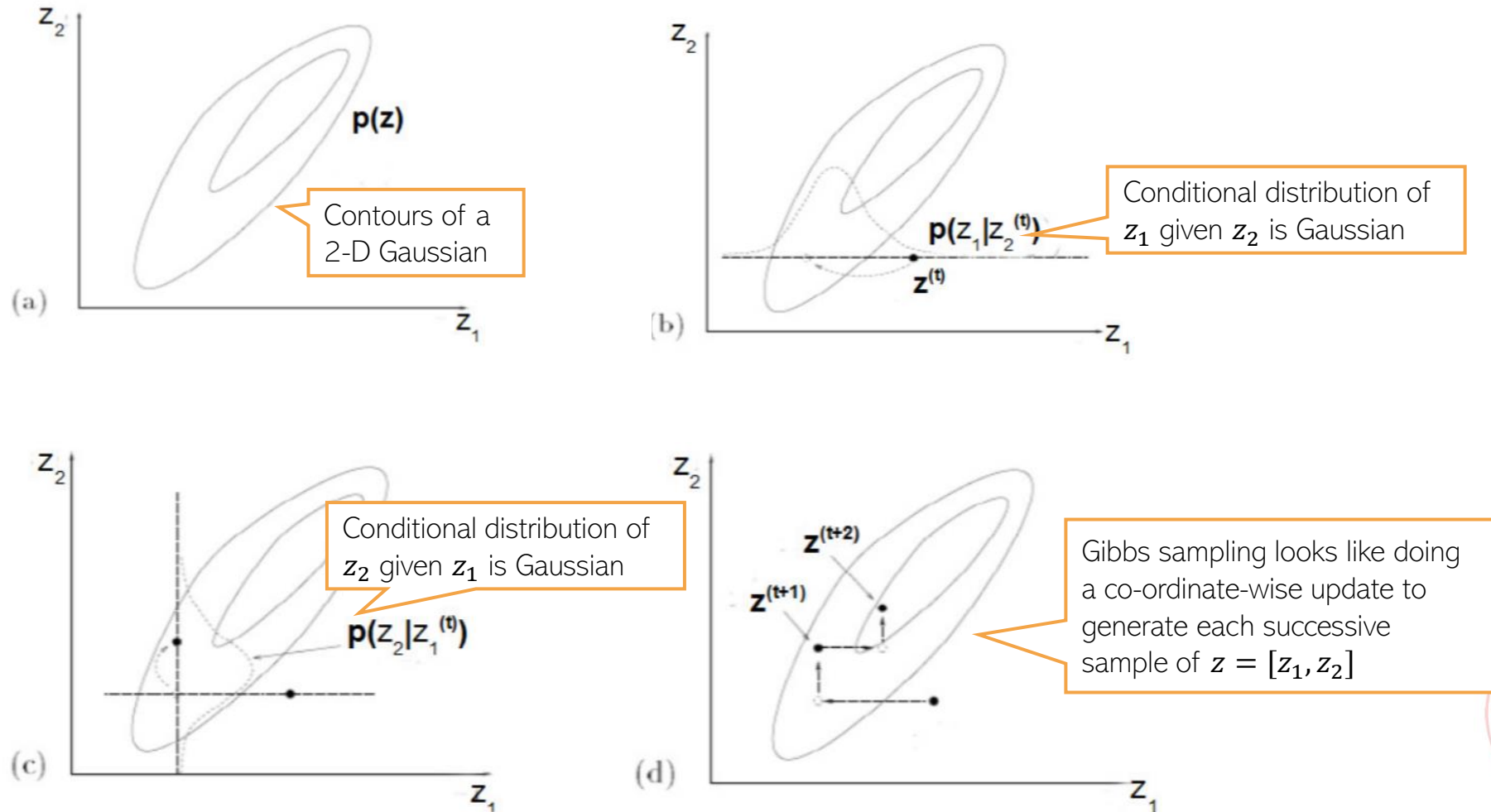
Each iteration will give us one sample $\mathbf{z}^{(\tau)}$ of $\mathbf{z} = [z_1, z_2, \dots, z_M]$

- 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



Gibbs Sampling: Another Simple Example

- Bayesian linear regression: $p(y_n|\mathbf{x}_n, \mathbf{w}, \beta) = \mathcal{N}(y_n|\mathbf{w}^\top \mathbf{x}_n, \beta^{-1})$, $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|0, \lambda^{-1}I)$, $p(\lambda) = \text{Gamma}(\lambda|a, b)$, $p(\beta) = \text{Gamma}(\beta|c, d)$. Gibbs sampler for $p(\mathbf{w}, \lambda, \beta|\mathbf{X}, \mathbf{y})$ will be
- Initialize λ, β as $\lambda^{(0)}, \beta^{(0)}$. For iteration $t = 1, 2, \dots, T$
 - Generate a random sample of \mathbf{w} by sampling from its CP as

$$\mathbf{w}^{(t)} \sim \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)}) \quad \text{where}$$

$$\boldsymbol{\Sigma}^{(t-1)} = (\beta^{(t-1)}\mathbf{X}^\top\mathbf{X} + \lambda^{(t-1)})^{-1}$$

$$\boldsymbol{\mu}^{(t-1)} = \left(\mathbf{X}^\top\mathbf{X} + \frac{\lambda^{(t-1)}}{\beta^{(t-1)}}\right)^{-1} \mathbf{X}^\top\mathbf{y}$$

- Generate a random sample of λ by sampling from its CP as

$$\lambda^{(t)} \sim \text{Gamma}\left(\lambda|a + \frac{D}{2}, b + \frac{\mathbf{w}^{(t)\top}\mathbf{w}^{(t)}}{2}\right)$$

- Generate a random sample of β by sampling from its CP as

$$\beta^{(t)} \sim \text{Gamma}\left(\beta|c + \frac{N}{2}, d + \frac{\|\mathbf{y} - \mathbf{X}\mathbf{w}^{(t)}\|^2}{2}\right)$$

- The posterior's approximation is the set of collected samples



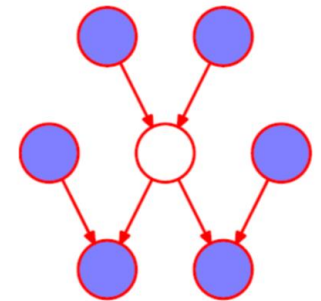
Gibbs Sampling: Some Comments

- One of the most popular MCMC algorithms
- Very easy to derive and implement for locally conjugate models
- Many variations exist, e.g.,
 - **Blocked Gibbs**: sample more than one component jointly (sometimes possible)
 - **Rao-Blackwellized Gibbs**: Can collapse (i.e., integrate out) the unneeded components while sampling. Also called “collapsed” Gibbs sampling
 - **MH within Gibbs**: If CPs are not easy to sample distributions
- Instead of sampling from CPs, an alternative is to use the mode of the CPs
 - Called the “**Iterative Conditional Mode**” (ICM) algorithm
 - ICM doesn't give the posterior though – it's more like ALT-OPT to get (approx) MAP estimate



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

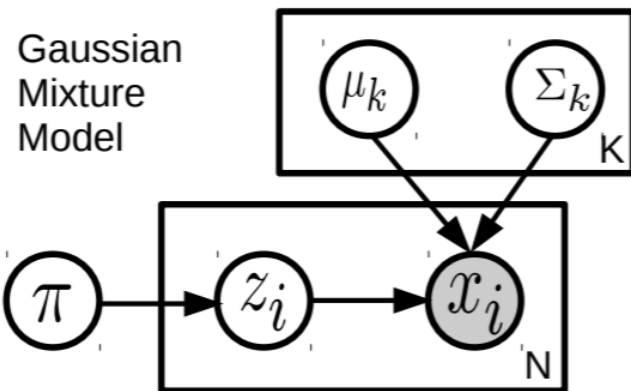


Markov Blanket



Gibbs Sampling: An Example

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



Joint distribution of data and unknowns

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

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

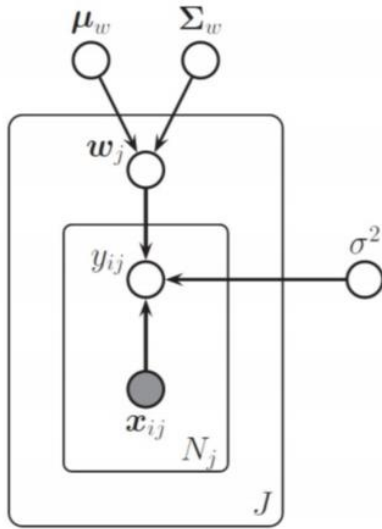
$$p(\boldsymbol{\pi} | \mathbf{z}) = \text{Dir}(\{\alpha_k + \sum_{i=1}^N \mathbb{I}(z_i = k)\}_{k=1}^K)$$

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

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

Gibbs Sampling: Another Example

J schools
Regression
Problem



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

Can verify that
Markov Blanket
property holds
for each CP

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

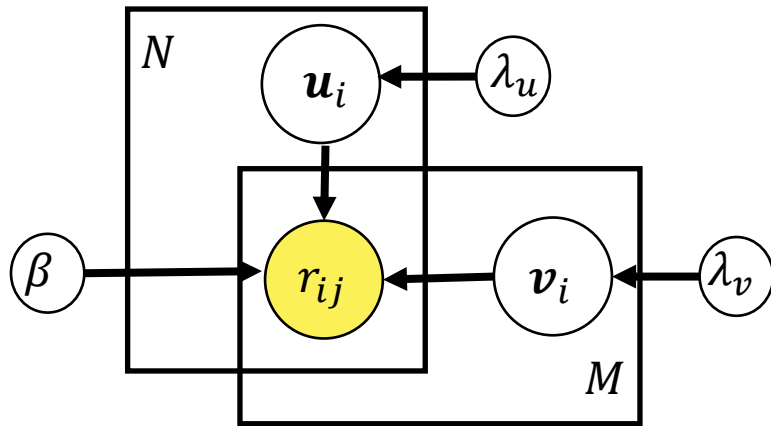
$$\begin{aligned}
 p(\boldsymbol{\mu}_w | \mathbf{w}_{1:J}, \boldsymbol{\Sigma}_w) &= \mathcal{N}(\boldsymbol{\mu} | \boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N) \\
 \boldsymbol{\Sigma}_N^{-1} &= \mathbf{V}_0^{-1} + J \boldsymbol{\Sigma}^{-1} \\
 \boldsymbol{\mu}_N &= \boldsymbol{\Sigma}_N (\mathbf{V}_0^{-1}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

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

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



Gibbs Sampling: One More 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

Using MCMC 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

