

Gibbs Sampling

initialize $z_1^0, \dots, z_n^0, \mu_{1:k}^0, \theta^0$

for iteration $m = 1, 2, \dots, M$:

$$z_{1:n}^m \sim p(z_{1:n} \mid \mu_{1:k}^{m-1}, \theta^{m-1}, x_{1:n})$$

$$\mu_{1:k}^m \sim p(\mu_{1:k} \mid z_{1:n}^m, x_{1:n})$$

$$\theta^m \sim p(\theta \mid z_{1:n}^m, x_{1:n})$$

(This should feel like EM.)

This returns a set of samples

$$\left\{ z_{1:n}^m, \mu_{1:k}^m, \theta^m \right\}_{m=1}^M$$

Claim:

$$\lim_{M \rightarrow \infty} \frac{1}{M} \sum_{m=1}^M \mathbb{1}(\mu_k^m \in A) = p(\mu_k \in A \mid x_{1:n})$$

for any subset A

Another way of saying this is that:

$$\lim_{M \rightarrow \infty} \Pr(\mu_k^M) = p(\mu_k \mid x_{1:n})$$

More generally:

$$\lim_{M \rightarrow \infty} \frac{1}{M} \sum_{m=1}^M f(z_{1:n}^m, \mu_{1:k}^m, \theta^m)$$

$$= \mathbb{E}_{\text{Posterior}} \left[f(z_{1:n}, \mu_{1:k}, \theta) \mid x_{1:n} \right]$$

Conceptually our algorithm is a Markov chain.

$$\text{State: } S^m \triangleq (z_{1:n}^m, \mu_{1:k}^m, \theta^m)$$

A Markov chain is a joint distribution of sequential random variables S_1, \dots, S_M which factorizes:

$$\pi(s^0, s^1, \dots, s^M) = \underbrace{\pi(s^0)}_{\text{"initial dist."}} \prod_{m=1}^M \underbrace{\pi(s^m | s^{m-1})}_{\text{"transition operator"}}$$

If $\pi(s^m = s | s^{m-1} = s')$ is the same for all m then the Markov chain is homogeneous. In this case:

$$\pi(s^m = s | s^{m-1} = s') \equiv \pi(s | s')$$

Each state $S^m \equiv (s_1^m, \dots, s_D^m)$ has D components.

$$\text{e.g. } S^m \equiv (z_1^m, \dots, z_n^m, \mu_1^m, \dots, \mu_k^m, \theta^m), \text{ so } D = n + k + 1$$

The m^{th} marginal of the chain is:

$$\begin{aligned} \pi(S^m = s) &= \int \pi(s^{m-1} = s') \pi(s^m = s | s^{m-1} = s') ds' \\ &\equiv \int \pi(s^{m-1} = s') \pi(s | s') ds' \end{aligned}$$

A distribution $\pi^*(s)$ is a stationary distribution of π if:

$$\pi^*(s) = \int \pi^*(s') \pi(s | s') ds' \equiv \mathbb{E}_{s' \sim \pi^*} [\pi(s | s')]$$

We also say that π^* is invariant to π .

Detailed balance

How can we know if π^* is a stationary dist.?

A sufficient condition is if it satisfies detailed balance:

$$\pi^*(s') \pi(s|s') = \pi^*(s) \pi(s'|s)$$

To see this condition implies that π^* is stationary:

$$\begin{aligned} \int \pi^*(s') \pi(s|s') ds' &= \int \pi^*(s) \pi(s'|s) ds' \\ &= \pi^*(s) \end{aligned}$$

Ergodicity

DB implies that π^* is a stationary dist. but not necessarily a unique one.

A Markov chain that is ergodic has a unique π^* .

$$\lim_{n \rightarrow \infty} \pi(s^n = s) = \pi^*(s), \text{ for any } s^0$$

A sufficient condition for ergodicity is that:

$$\pi(s|s') > 0 \quad \forall s, s'$$

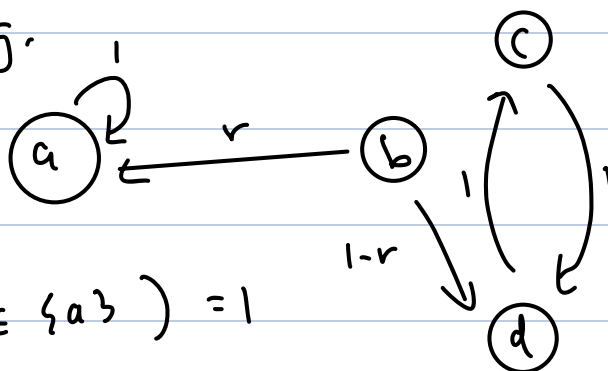
More generally, a chain is ergodic if it is

① Irreducible

② Aperiodic

Irreducibility

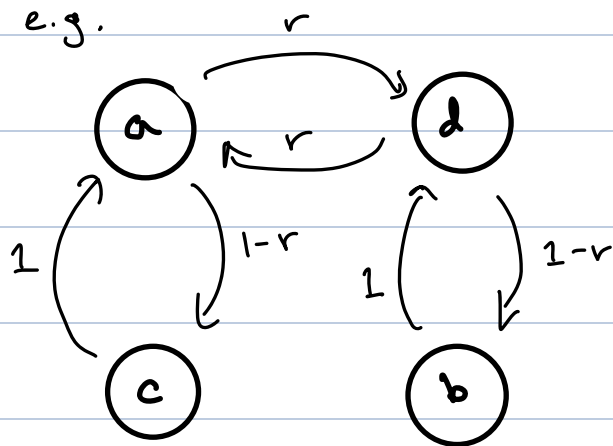
e.g.



- $P(S^m \in \{a\} \mid S^{m-1} \in \{a\}) = 1$
- $P(S^m \in \{c, d\} \mid S^{m-1} \in \{c, d\}) = 1$
- $\{a\}$ and $\{c, d\}$ are recurrent classes because all nodes are reachable from all other nodes within the class
- $\{a, b, c, d\}$ is not a recurrent class, since e.g.
 $P(S^m = b \mid S^0 = c) = 0 \quad \forall m \geq 1$

Periodicity

e.g.



- $P(S^m \in \{c, d\} \mid S^{m-1} \in \{a, b\}) = 1$
- $P(S^{m+1} \in \{a, b\} \mid S^m \in \{c, d\}) = 1$
- So we can partition the chain into two periodic classes $\{a, b\}$, $\{c, d\}$.

The Gibbs sampler is a Markov chain with transitions:

$$\begin{aligned}\pi_{\text{Gibbs}}(s^m | s^{m-1}) &= p(s_1^m | s_2^{m-1}, \dots, s_D^{m-1}, x) \\ &\quad p(s_2^m | s_1^m, s_3^{m-1}, \dots, s_D^{m-1}, x) \\ &\quad \vdots \\ &\quad p(s_D^m | s_1^m, \dots, s_{D-1}^m, x)\end{aligned}$$

where $p(s_d | s_1, \dots, s_{d-1}, s_{d+1}, \dots, s_D, x)$ is the complete conditional of latent variable s_d .

(These conditionals also condition on the data x)

For any coordinate d , WLOG, say that it is sampled first:

$$\pi_{\text{Gibbs}}(s_d^m = s_d) = \int p(s_d | s'_{-d}, \bar{x}) \pi_{\text{Gibbs}}^{m-1}(s'_{-d} = s'_{-d}) ds'_{-d}$$

$$\text{say that } \pi_{\text{Gibbs}}^{m-1}(s'_{-d} = s'_{-d}) = p(s'_{-d} | x)$$

$$= p(s_d | \bar{x})$$

So if the marginal for s_{-d}^{m-1} is the posterior then the marginal for s_d^m is too.

It's easier to see the implication for $D=2$:

$$\begin{aligned}\pi_{\text{Gibbs}}(s^m = (s_1, s_2)) &= \int p(s_1 | s_2, x) p(s_2 | s_1', x) \pi_{\text{Gibbs}}^{m-1}(s_1', s_2') ds_1' ds_2' \\ &= p(s_1 | s_2, x) \int p(s_2 | s_1', x) \pi_{\text{Gibbs}}^{m-1}(s_1' = s_1) ds_1\end{aligned}$$

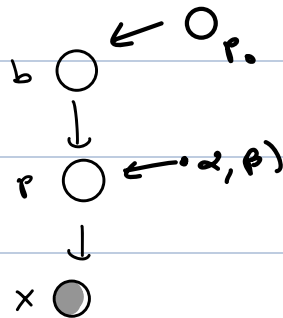
$$\begin{aligned}\text{If } \pi_{\text{Gibbs}}(s^{m-1} = (s_1', s_2')) &= p(s_1', s_2' | x) \\ &= p(s_1, s_2 | x)\end{aligned}$$

If the Gibbs chain is ergodic, then:

$$\lim_{m \rightarrow \infty} \pi_{\text{Gibbs}}(s^m = s) = p(s | \bar{x})$$

An example of a non-ergodic Gibbs chain:

$$\begin{aligned}b &\sim \text{Bern}(p_0) \\ p &\sim \begin{cases} \delta_0 & \text{if } b=0 \\ \text{Beta}(\alpha, \beta) & \text{if } b=1 \end{cases} \\ x &\sim \begin{cases} \delta_0 & \text{if } p=0 \\ \text{Bern}(p) & \text{if } p>0 \end{cases}\end{aligned}$$



$$p(p=0 | b=0, \bar{x}=0) = 1$$

$$p(b=0 | p=0, \bar{x}=0) = 1$$

So this Gibbs sampler will never leave the state $(p=0, b=0)$ if $\bar{x}=0$.

Blocked and collapsed Gibbs

Blocking means we sample some latent variables from their joint conditional, e.g.,

$$\begin{aligned} \pi_{\text{Gibbs}}(p, b, p_0 \mid p', b', p_0') & \quad \swarrow (p, b) \text{ are } \underline{\text{blocked}} \\ &= P(p_0 \mid p, b, \bar{x}) P(p, b \mid p_0', \bar{x}) \end{aligned}$$

Related to this is collapsing where we sample a variable with another variable marginalized out

$$= P(p_0 \mid p, b, \bar{x}) P(p \mid b, p_0', \bar{x}) P(b \mid p_0', \bar{x})$$

\uparrow
 p is "collapsed" out

This chain is no longer non-ergodic.

Missing data

$$X = X_{\text{obs}} \cup X_{\text{miss}}$$

Treating X_{miss} as a latent variable

Gibbs w/ missing data

init $z^0, \theta^0, \Phi^0, X_{\text{miss}}^0$

for iteration m :

$$z^m \sim P(z \mid \theta^{m-1}, \Phi^{m-1}, X_{\text{miss}}^{m-1}, X_{\text{obs}})$$

$$\theta^m \sim \dots$$

$$\Phi^m \sim \dots$$

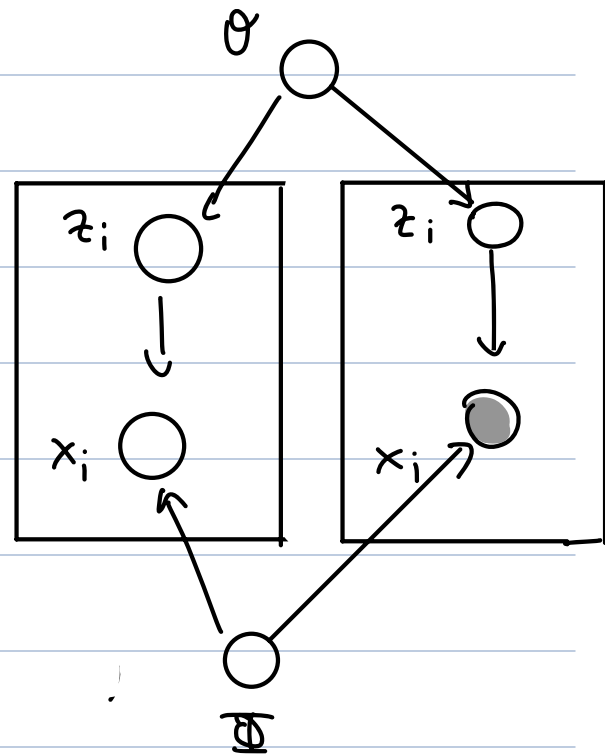
$$X_{\text{miss}}^m \sim P(X_{\text{miss}} \mid z^m, \theta^m, \Phi^m, X_{\text{obs}})$$

Question: how do we sample θ 's?

Question: $\lim_{m \rightarrow \infty} \frac{1}{m} \sum X_{\text{miss}}^m = ?$

$$= \mathbb{E}[X_{\text{miss}} \mid X_{\text{obs}}]$$

(posterior predictive expectation)



Geweke testing (2004)

Consider $X = X_{\text{miss}}$ (X_{obs} is empty).

What would be the stationary dist. of the Gibbs sampler above?

$$S \triangleq \{z, \emptyset, \oplus\}$$

$$\pi^*(S, X_{\text{miss}}) = P(S, X_{\text{miss}} \mid X_{\text{obs}}) \equiv P(S, X)$$

"X" "empty" ↑ the joint!

This suggests a way to test our Gibbs sampler.

Define the forward sampler to be:

for $m = 1 \dots M$:

$$S_f^m \sim P(S)$$

$$X_f^m \sim P(X \mid S_f^m)$$

And the backward sampler to be Gibbs

for $m = 1 \dots M$:

$$S_b^m \sim \pi(S_b^m \mid S_b^{m-1}, X_b^{m-1})$$

$$X_b^m \sim P(X \mid S_b^m)$$

These should both produce samples from the joint:

$$(X_f^m, S_f^m) \sim P(X, S)$$

$$(X_b^m, S_b^m) \sim P(X, S)$$

We can implement both and test whether the samples they generate have the same distribution. If not, there's a bug.

Collapsed Gibbs in the expfam mixture

$$P(z_i = k) = \vartheta_k$$

$$P(x_i | z_i = j) = G(x_i; \eta_j)$$

$$F(\eta_k; \lambda)$$

Collapse out η_1, \dots, η_k when Gibbs sampling z_1, \dots, z_n :

$$P(z_i = j | x, z_{-i}, \theta)$$

$$\propto \vartheta_j P(x_i | z_i = j, z_{-i}, x_{-i})$$

$$\propto \vartheta_j \int P(x_i, \eta_j | z_i = j, z_{-i}, x_{-i}) d\eta_j$$

$$\propto \vartheta_j \underbrace{\int P(x_i | \eta_j) P(\eta_j | z_{-i}, x_{-i}) d\eta_j}_{\text{posterior prediction}} = G(x_i; \eta_j) F(\eta_j; \lambda_{j,-i}) d\eta_j$$

$$\lambda_{j,-i} \equiv \begin{bmatrix} \lambda_{j,-i,1} \\ \lambda_{j,-i,2} \end{bmatrix} \stackrel{\triangle}{=} \begin{bmatrix} \lambda_1 + \sum_{i' \neq i} t(x_{i'}) f_{i,j} \\ \lambda_2 + \sum_{i' \neq i} f_{i,j} \end{bmatrix}$$

$$\propto \vartheta_j \frac{\exp(a_e([\lambda_{j,-i,1} + t(x_i), \lambda_{j,-i,2} + 1]))}{\exp(a_e([\lambda_{j,-i,1}, \lambda_{j,-i,2}]))}$$

