

Hierarchical model for "8 schools" (from last time)

The generative process:

$$\tau^2 \sim \chi^2(\nu_0, \tau_0^2)$$

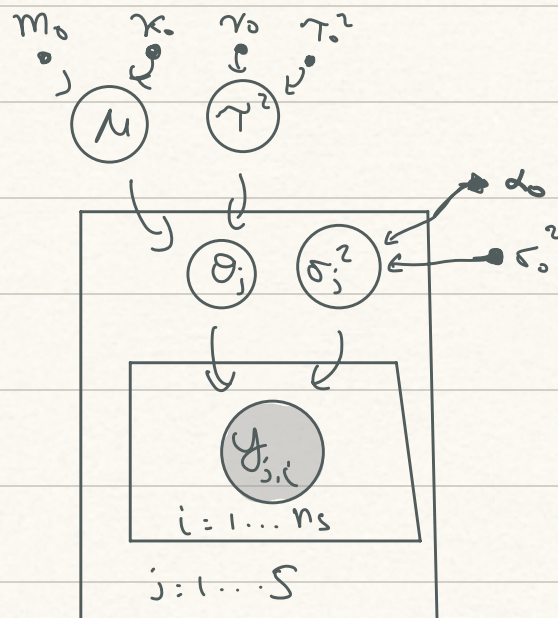
$$\mu \sim \mathcal{N}(m_0, 1/\kappa_0)$$

$$\sigma_j^2 \stackrel{iid}{\sim} \chi^2(\alpha_0, \sigma_0^2) \quad j=1 \dots S$$

$$\theta_j \stackrel{iid}{\sim} \mathcal{N}(\mu, \tau^2) \quad j=1 \dots S$$

$$y_{j,i} \stackrel{iid}{\sim} \mathcal{N}(\theta_j, \sigma_j^2) \quad i=1 \dots n_s$$

The probabilistic graphical model (PGM):



The set of latent variables is:

$$\mathcal{Z} = \{ (\theta_j, \sigma_j^2)_{j=1}^S, \mu, \tau^2 \}, \quad |\mathcal{Z}| = D (= 2S + 2, \text{ in this case})$$

The set of hyperparameters is:

$$\eta_0 = \{ m_0, \kappa_0, \nu_0, \tau_0^2, \alpha_0, \sigma_0^2 \}$$

The joint posterior in this model is intractable

...but complete conditionals are tractable

$$P(\mathcal{Z} | \mathcal{Y}, \eta_0) \quad \times$$

$$P(\mathcal{Z}_d | \mathcal{Z}_{\setminus d}, \mathcal{Y}, \eta_0) \quad \checkmark$$

For example, by Gaussian-Gaussian conjugacy:

$$P(\mu | -) = \mathcal{N}(m_S, 1/\kappa_S)$$

$$m_S = w_S \frac{1}{S} \sum_{j=1}^S \theta_j + (1 - w_S) m_0$$

$$w_S = \frac{S/\tau^2}{\kappa_0 + S/\tau^2}$$

$$\kappa_S = \kappa_0 + S/\tau^2$$

The Gibbs sampler (today):

for $m = 1 \dots M$:

for $d = 1 \dots D$:

$$\mathcal{Z}_d^m \sim P(\mathcal{Z}_d | -)$$

Claim:

$$\lim_{M \rightarrow \infty} \frac{1}{M} \sum_{m=1}^M f(\mathcal{Z}^m) = \mathbb{E}_{P(\mathcal{Z} | \mathcal{Y})} [f(\mathcal{Z})]$$

PGM semantics

Notice that the complete conditional for μ depends on τ^2 , θ_j , but not Y . Why?

$$\begin{aligned}
 P(\mu | -) &= P(\mu | z_\mu, Y, \eta) = \frac{P(\mu, z_\mu, Y | \eta)}{\int P(\mu, z_\mu, Y | \eta) d\mu} \\
 &= \frac{P(\mu) P(\tau^2) \left[\prod_j P(\delta_j | \mu, \tau^2) P(\delta_j^2) \right] \prod_{i,j} P(y_{ij} | \theta_j, \delta_j^2)}{\int P(\mu) P(\tau^2) \left[\prod_j P(\delta_j | \mu, \tau^2) P(\delta_j^2) \right] \prod_{i,j} P(y_{ij} | \theta_j, \delta_j^2) d\mu} \\
 &= \frac{P(\mu) P(\tau^2)}{\int P(\mu) P(\tau^2) d\mu} = \int P(\mu) P(\tau^2) d\mu
 \end{aligned}$$

The kernel of the posterior involves τ^2 and θ_j but not Y ; this is due to the conditional independence assumptions of the model.

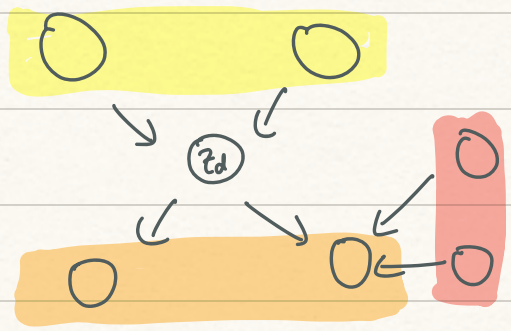
e.g. $y_{ij} \perp \mu \mid \theta_j, \delta_j^2$

PGMs encode the conditional independences in a model. In general, the joint distribution is:

$$\begin{aligned}
 P(z, y) &= \left[\prod_{d=1}^D P(z_d | z_{\text{par}(z_d)}) \right] P(y | z_{\text{par}(y)}) \\
 \text{par}(z_d) &= \{ z_{d'} : z_{d'} \rightarrow z_d \in \text{graph} \} \\
 &\quad \text{"parents" of } z_d
 \end{aligned}$$

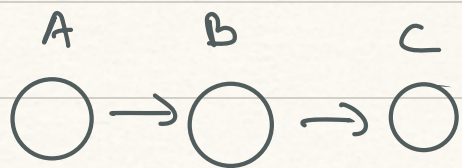
A variable is conditionally independent of all other variables given its **Markov blanket**:

$$\text{MB}(z_d) = \{ \underbrace{\text{par}(z_d)}_{\text{"parents"}} \cup \underbrace{\text{child}(z_d)}_{\text{"children"}} \cup \underbrace{\text{coparents}(z_d)}_{\text{"co-parents"}} \}$$

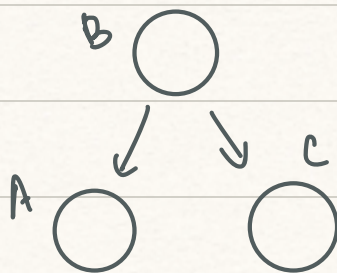


$$\begin{aligned}
 z_d &\perp z_{d'} \mid \text{MB}(z_d) \\
 \forall z_{d'} &\notin \text{MB}(z_d)
 \end{aligned}$$

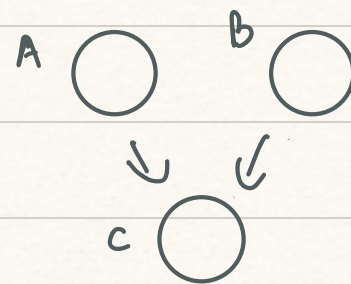
Directional separation (D-separation)



① chain



② tree



③ V-structure

$$A \perp\!\!\!\perp B \mid C$$

$$A \not\perp\!\!\!\perp B$$

$$A \perp\!\!\!\perp C \mid B$$

$$A \not\perp\!\!\!\perp C$$

$$A \not\perp\!\!\!\perp B \mid C$$

$$A \perp\!\!\!\perp B$$

There are 3 basic "motifs" (i.e., mini-graphs) that encode fundamentally different conditional independence semantics.

1) Chain: e.g., "the past is independent of the future given the present".

2) Tree: e.g., iid sampling: student 1 independent of student 2 given school mean/variance

3) V-structure: "explaining away", e.g., COVID not independent of cold given coughing

We will talk more about PGM semantics in week 5.

Monte carlo approximations

Last time, we motivated the posterior expectation of the school means:

$$\hat{\theta} = \mathbb{E}_{p(\theta|y)} [\theta]$$

where the expectation is with respect to the posterior marginal:

$$p(\theta | y) = \int p(\theta, z_{- \theta} | y) dz_{- \theta}$$

↑ posterior marginal
marginalizes out all other latent

In complicated enough models, we cannot analytically compute these marginals.

However, we can still approximate posterior expectations using Monte Carlo.

$$\mathbb{E}_{p(\theta|y)} [f(\theta)] = \frac{1}{M} \sum_{m=1}^M f(\theta^m),$$

$\theta^m \stackrel{iid}{\sim} p(\theta | y), m=1 \dots M$

Clearly the Monte Carlo estimator is unbiased:

$$\mathbb{E}_{p(\theta|y)} \left[\frac{1}{M} \sum_{m=1}^M f(\theta^m) \right] = \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{p(\theta|y)} [f(\theta)] = \mathbb{E} [f(\theta)]$$

It has variance:

$$\text{Var} \left[\frac{1}{M} \sum_{j=1}^M f(\theta_j^m) \right] = \frac{1}{M^2} \left(\sum_m \text{Var}(f(\theta)) + 2 \sum_{m < m'} \text{Cov}[f(\theta^m), f(\theta^{m'})] \right)$$

If the samples are uncorrelated, the root mean squared error (RMSE) of the Monte Carlo estimator is then:

$$\text{RMSE} \left(\frac{1}{M} \sum_{m=1}^M f(\theta^m) \right) = \sqrt{\frac{1}{M} \text{Var}(f(\theta))} = O\left(\frac{1}{\sqrt{M}}\right)$$

Notice that this does not depend on the dimension of the parameter.

Takeaway: if we can draw (uncorrected) iid samples from the posterior, we can get unbiased estimates of posterior expectations whose error does not depend on the dimension of our parameter set. This is the basis for MCMC methods.

Markov chains and MCMC:

The Gibbs sampler is a **Markov chain** with state equal to the latent variables:

$$\text{State: } S^m = (S_1^m \dots S_D^m)$$

A Markov chain defines a joint distribution over a sequence of states, which factorizes:

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

A Markov chain is "**homogenous**" if the transition operator is the same at all steps:

$$\pi(S^m = s | S^{m-1} = s') \equiv \pi(s | s') \text{ for "homogenous" chains}$$

The marginal distribution over the state at step m is:

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

A "**stationary distribution**" of the Markov chain is one with the following property:

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

A stationary distribution is "**invariant**" to the transition operator.

More informally, if you start the chain from a stationary distribution, you stay there.

A Gibbs sampler is one kind of **Markov Chain Monte Carlo (MCMC)** method.

MCMC rely on **Markov chains** whose stationary distribution is the exact posterior.

How do we know if a distribution is a stationary distribution of a Markov chain?

A sufficient condition is "**detailed balance**":

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

To see why, we integrate both sides:

$$\begin{aligned} \underbrace{\int \pi^*(s') \pi(s | s') ds'}_{= \pi^*(s) \text{ (by definition)}} &= \underbrace{\int \pi^*(s) \pi(s' | s) ds'}_{= \pi^*(s) \int \pi(s' | s) ds'} \end{aligned}$$

The target distribution is thus

an stationary distribution may have multiple stationary distributions.

A Markov chain has a unique stationary distribution if it is "ergodic", and if it is ergodic, it will eventually converge to its unique stationary distribution:

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

A sufficient condition for ergodicity is that all states are reachable from all other states under the transition operator:

$$\pi(s | s') > 0$$

More generally, a chain is ergodic if it is 1) irreducible and 2) aperiodic.
(We will not define these in lecture, unless there is time.)

Gibbs sampler:

A Gibbs sampler is a Markov chain whose transition operator updates one coordinate at a time from the conditional distribution of the target distribution.

$$\pi_{\text{Gibbs}}(S^m = s \mid S^{m-1} = s') = \prod_{d=1}^D P^*(S_d = s_d \mid S_{\setminus d} = s_{\setminus d}, S_{>d} = s'_{>d})$$

The marginal distribution over the d th coordinate (i.e., latent variable) is:

$$\pi_{\text{Gibbs}}(S_d^m = s_d) = \int P^*(S_d \mid s'_{\setminus d}) \pi_{\text{Gibbs}}(S^{m-1} = s') ds'$$

Say that at $(m-1)$ th marginal is the target distribution:

$$\text{Say } \pi_{\text{Gibbs}}(S^{m-1} = s') = P^*(s'_1 \dots s'_D).$$

Then the m th marginal over the d th coordinate is the marginal under the target:

$$\begin{aligned} &= \int ds' P^*(s_d \mid s'_{\setminus d}) P^*(s'_1 \dots s'_D) \\ &= \int ds'_{\setminus d} P^*(s_d \mid s'_{\setminus d}) \int ds'_d P^*(s'_1 \dots s'_D) \\ &= P^*(s_d) \end{aligned}$$

So : $P^*(s_1 \dots s_D)$ is a stationary dist.

In Gibbs sampling, our target distribution is the posterior, and our conditionals are the complete conditionals. If our chain is ergodic then, the posterior is the unique stationary distribution, and the chain will eventually converge to it.

$$\lim_{m \rightarrow \infty} \prod_{\text{Gibbs}} (S^m = s) = P^*(s) = P(Z|Y)$$

Non-ergodic chains usually arise when there are deterministic steps in the generative process. Consider the following model:

$$b \sim \text{Bernoulli}(p)$$

Consider the following complete conditionals for p and b given $y=0$:

$$p \sim \begin{cases} \delta_0 & \text{if } b=0 \\ \text{Beta}(\alpha, \eta) & \text{if } b=1 \end{cases}$$

$$P(p=0 \mid b=0, y=0) = 1$$

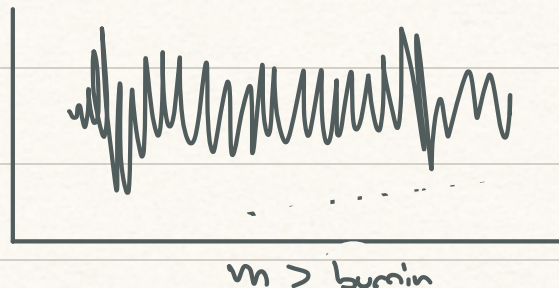
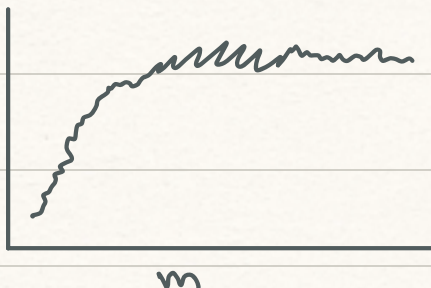
$$P(b=0 \mid p=0, y=0) = 1$$

$$y \sim \begin{cases} \delta_0 & \text{if } p=0 \\ \text{Bernoulli}(p) & \text{if } p>0 \end{cases}$$

So if the chain starts at $(b=0, p=0)$, it will never leave, and it won't converge to the posterior. As a rule of thumb, if all complete conditionals in your model always place non-zero probability on all possible outcomes, then the Gibbs sampler will be ergodic.

Practical considerations:

The chain will eventually be independent of the initial state, but that is only after enough transitions. In practice, we set a "burn-in" or "warm-up" period where we run the chain until it appears to converge; then we start collecting samples.



"trace plots"

As mentioned above, our Monte Carlo estimates will have higher error if samples are correlated. Samples in a Gibbs chain are indeed correlated, since they are directly conditioned on each other. To reduce our chain's "autocorrelation" we will "thin" the chain by retaining only every H samples (e.g., every 50th sample).

Missing data:

What if we have missing / heldout data?

$$Y = Y_{\text{obs}} \cup Y_{\text{miss}}, \quad S = \{z, y_{\text{miss}}\}$$

MCMC has a natural way to handle this: by simply treating missing observations like any of the other latent variables. The Gibbs sampler then becomes:

for $m = 1 \dots M$:

for $d = 1 \dots D$:

$$z_d^m \sim P(z_d \mid z_{\setminus d}^m, z_{\setminus d}^{m-1}, y_{\text{miss}}^{m-1}, y_{\text{obs}})$$

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

This last step samples the missing data from its complete conditional.

Question: in the "8 schools" model, what is this conditional?

Question: What expectations do the following Monte Carlo estimates approximate:

$$\lim_{M \rightarrow \infty} \frac{1}{M} \sum_{m=1}^M y_{\text{miss}}^m = ?$$

$$\lim_{M \rightarrow \infty} \frac{1}{M} \sum_{m=1}^M P(y_{\text{miss}} \mid z^m, y_{\text{obs}}) = ?$$

Question: What would the Gibbs stationary distribution be in this case?

$$\pi_{\text{Gibbs}}^{\infty}(S) = ?$$

Testing MCMC:

What if we treated all our data as missing? This is strange but worthwhile:

$$Y = Y_{\text{miss}}, \quad S = \{Y, Z\}$$

What would the Gibbs stationary distribution be in this case?

$$\pi_{\text{Gibbs}}^*(S) = P(Y_{\text{miss}}, Z \mid Y_{\text{obs}}) = P(Y, Z)$$

full joint! ↗

This observation is the basis for a very useful trick called “Geweke testing” which lets you test whether your implementation of an MCMC algorithm is correct.

Define two different sampling algorithms called the **forward** and **backward** sampler.

Forward sampler

Backward sampler

for $m=1 \dots M$:

for $m=1 \dots M$:

$$Z_f^m \sim P(Z)$$

$$Z_b^m \sim \pi_{\text{Gibbs}}(Z \mid Z_b^{m-1}, Y_b^{m-1})$$

$$Y_f^m \sim P(Y \mid Z)$$

$$Y_b^m \sim P(Y \mid Z_b^m)$$



$$(Z_f^m, Y_f^m) \sim P(Z, Y)$$

$$(Z_b^m, Y_b^m) \sim P(Z, Y)$$

The key idea is that both should be drawing samples from the joint. The forward sampler simply samples M independent joint samples by sampling from the prior and then the likelihood. The backward sampler begins at some initialization, then it iteratively resamples the latents given data and data given latents. The backward sampler involves the Gibbs transition operator while the forward does not. Thus you can test whether your transition operator is correct by comparing the joint samples from each (they should be drawn from the same distribution!).

