

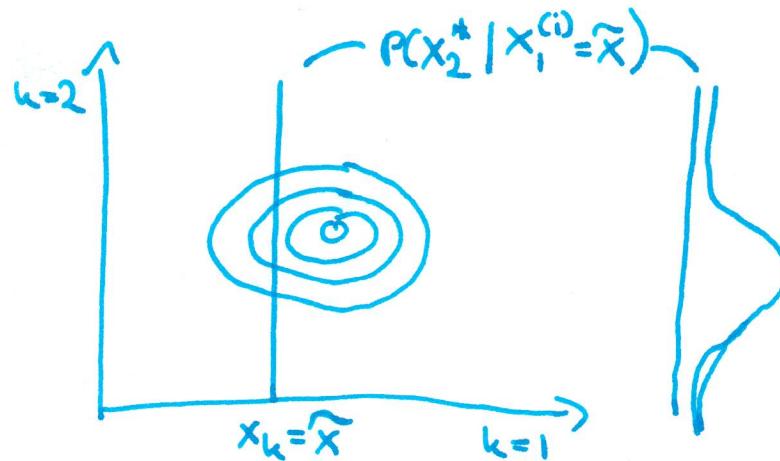
Lecture 10 - More Gibbs Sampling

Let's prove that Gibbs sampling converges to our target distribution.

Notation : $x_{-k}^{(i)}$ all components of $x^{(i)}$ except k
 $x_k^{(i)}$ k^{th} component of $x^{(i)}$

We said that Gibbs sampling is a form of Metropolis-Hastings where we always accept the proposed new position, and this new position is sampled from the conditional:

$$q_k(x^* | x^{(i)}) = \begin{cases} p(x_k^* | x_{-k}^{(i)}) & \text{for } x_k^* = x_{-k}^{(i)}, \\ 0 & \text{otherwise} \end{cases}$$



For the reverse direction we have:

$$q_u(x^{(i)} | x^*) = \begin{cases} p(x_u^{(i)} | x_{-k}^*) & \text{for } x_{-k}^{(i)} = x_{-k}^* \\ 0 & \text{otherwise} \end{cases}$$

Let's look at our M-H acceptance probability:

$$A = \min \left\{ 1, \frac{p(x^*)}{p(x^{(i)})} \cdot \frac{q_u(x^{(i)} | x^*)}{q_u(x^* | x^{(i)})} \right\}$$

$$\begin{aligned} \text{we know that: } p(x^*) &= p(x_{-k}^*, x_k^*) \\ &= p(x_{-k}^* | x_k^*) \cdot p(x_k^*) \\ &= p(x_k^* | x_{-k}^*) \cdot p(x_{-k}^*) \end{aligned}$$

$$\Rightarrow A = \min \left\{ 1, \frac{p(x_k^* | x_{-k}^*) \cdot p(x_{-k}^*) \cdot q_u(x^{(i)} | x^*)}{p(x_u^{(i)} | x_{-k}^{(i)}) \cdot p(x_{-k}^{(i)}) \cdot q_u(x^* | x^{(i)})} \right\}$$

Now we put in the definition of
 $q_u(x^{(i)} | x^*)$ and $q_u(x^* | x^{(i)})$

$$A = \min \left\{ 1, \frac{p(x_k^* | x_{-k}^*) \cdot p(x_{-k}^*) \cdot p(x_k^{(i)} | x_{-k}^*)}{p(x_k^{(i)} | x_{-k}^{(i)}) \cdot p(x_{-k}^{(i)}) \cdot p(x_k^* | x_{-k}^{(i)})} \right\}$$

Next we realize that $x_{-k}^{(i)} = x_{-k}^*$ because the new position for x^* is computed componentwise, so x^* differs from $x^{(i)}$ only in the new k th component.

$$\frac{\cancel{p(x_k^* | x_{-k}^{(i)})} \cdot \cancel{p(x_{-k}^{(i)})} \cdot \cancel{p(x_k^{(i)} | x_{-k}^*)}}{\cancel{p(x_k^{(i)} | x_{-k}^{(i)})} \cdot \cancel{p(x_{-k}^{(i)})} \cdot \cancel{p(x_k^* | x_{-k}^{(i)})}} = 1$$

$$\Rightarrow A = \min \{1, 1\} = 1$$

So we can rest assured that Gibbs sampling is a form of Metropolis-Hastings with acceptance probability 1.

Rat tumor example

- We are concerned about the tumor rates of rats in control group studies.
- Θ is the probability that a rat receiving no treatment develops a tumor.
- Data: $n=14$ rats $y=4$ rats develop tumor
 \Rightarrow We can solve this analog to the infected people in a city example.

Naive estimate: $\frac{4}{14} = 0.286$

Bayes estimate:

Likelihood: $P(Y_i|\Theta_i, n_i) = \text{Binomial}(n_i, y_i, \Theta_i)$
 $= \binom{n_i}{y_i} \cdot \Theta_i^{y_i} \cdot (1-\Theta_i)^{n_i-y_i}$

Conjugate prior: $P(\Theta_i) = \text{Beta}(\Theta_i, \alpha, \beta)$

$$= \frac{\Theta_i^{\alpha-1} (1-\Theta_i)^{\beta-1}}{B(\alpha, \beta)}$$

\hookrightarrow Beta function

- To solve this we need to choose α and β
- ⇒ In the beginning of the lecture we needed an expert to give us values.
- ⇒ This time we look for more data instead.
- ⇒ We find 10 studies in the literature

Now our data looks like this:

$$\begin{aligned} n &: [n_1, n_2, \dots, n_{10}] \\ y &: [y_1, y_2, \dots, y_{10}] \end{aligned} \Rightarrow \text{we are saving our new data point for next week.}$$

What can we do with this data?

$$\text{Likelihood: } P(y|\theta, n) = \prod_i P(y_i|\theta_i, n_i)$$

We need to think a bit more about θ_0 , α and β .

We have three choices:

- complete pooling
- unpooling
- partial pooling

Complete pooling:

- We assume there is only one tumor rate and $\theta_i = \hat{\theta}$, $\alpha_i = \hat{\alpha}$, $\beta_i = \hat{\beta}$.
- This is likely to underfit the data, as we have experiments from different labs, times, and rat lineages.

Unpooling:

- We are completely flexible, allowing different θ_i , α_i and β_i for each experiment:
$$\theta_i \sim \text{Beta}(\alpha_i, \beta_i)$$
- These are a lot of parameters!
- We are likely to overfit with this model.

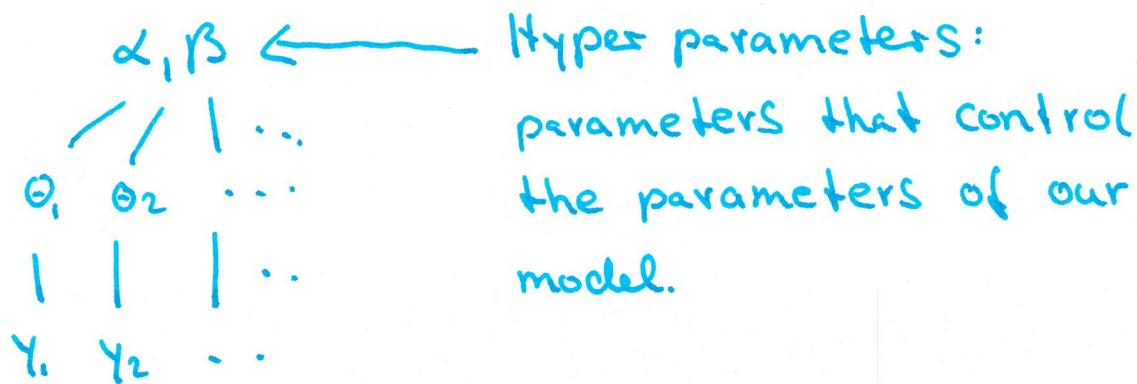
Partial pooling:

- While conditions vary a bit, all experiments still come from labs with similar conditions, and all experiments are from rats.

\Rightarrow We allow different Θ_i for each experiment, but we assume they all come from the same distribution:

$$\Theta_{gi} \sim \text{Beta}(\alpha, \beta)$$

\Rightarrow This is a hierarchical Bayes model. The Bioassay was a simplified version of this where the relationship between Θ and α, β was deterministic.



\Rightarrow We now need to find 10 tumor rates Θ_i and α, β

\Rightarrow We have a 72 dimensional problem.

Our complete posterior:

$$P(\theta, \alpha, \beta | Y, n) \propto \underbrace{P(Y | \theta, \alpha, \beta, n)}_{\text{Likelihood}} \cdot \underbrace{P(\theta, \alpha, \beta)}_{\text{Prior}}$$

We already have our likelihood as a Binomial
so we are back to thinking about the prior

Prior:

$$P(\theta, \alpha, \beta) = P(\alpha, \beta) \cdot \prod_i P(\theta_i | \alpha, \beta) \\ \hookrightarrow \text{Beta}(\theta_i; \alpha, \beta)$$

What about $p(\alpha, \beta)$?

- In the past we often just made it uniform
- In hierarchical models this can be tricky
- For large values of α and β our posterior goes to infinity and becomes non-integrable.

Let's think about α and β in a more intuitive way:

The mean of the Beta distribution is:

$$\mu = \frac{\alpha}{\alpha + \beta}$$

- It makes sense to have a uniform prior on the mean tumor rate, which corresponds to μ .
- In addition we put a uniform prior on $v = (\alpha + \beta)^{\frac{1}{2}}$ where $\alpha + \beta$ can be understood as our sample size.
- Note: to see this look at $\alpha + \beta$ for the Beta and the result when using it as a conjugate prior for our Binomial!

Before	After
$\text{Beta}(\alpha, \beta)$	$\text{Beta}(\alpha + y, \beta + N - y)$
$\alpha + \beta$	$\alpha + y + \beta + N - y = \alpha + \beta + N$

- What do these uniform priors on μ and v mean for α and β ?

$$\mu = \frac{\alpha}{\alpha + \beta} \quad v = \frac{1}{\sqrt{\alpha + \beta}} \quad v^2 = \frac{1}{\alpha + \beta}$$

For α we have:

$$\mu = \frac{\alpha}{\alpha + \beta}$$

$$\Leftrightarrow \frac{\mu}{\alpha} = \frac{1}{\alpha + \beta}$$

$$\textcircled{2} \quad \frac{\mu}{\lambda} = v^2$$

$$\Leftrightarrow \lambda = \frac{\mu}{v^2}$$

For β we have:

$$\frac{1}{v^2} = \lambda + \beta \Leftrightarrow \beta = \frac{1}{v^2} - \lambda \Leftrightarrow \beta = \frac{1}{v^2} - \frac{\mu}{v^2}$$

$$\Leftrightarrow \beta = \frac{1-\mu}{v^2}$$

So we got our inverse transform now, the mapping from $\mu, v \rightarrow \lambda, \beta$.

We need the Jacobian to get the factor by which this mapping expands or shrinks volumes.

Note: Wikipedia has really great articles for non-physicists who would like to understand this in more depth.

The Jacobian matrix for the inverse transform is:

$$\begin{vmatrix} \frac{\partial \lambda}{\partial \mu} & \frac{\partial \lambda}{\partial v} \\ \frac{\partial \beta}{\partial \mu} & \frac{\partial \beta}{\partial v} \end{vmatrix} = \begin{vmatrix} \frac{1}{v^2} & -2 \frac{\mu}{v^3} \\ -\frac{1}{v^2} & \frac{2(\mu-1)}{v^3} \end{vmatrix}$$

For the determinant we get:

$$\frac{\partial \mathcal{L}}{\partial \mu} \cdot \frac{\partial \beta}{\partial v} - \frac{\partial \mathcal{L}}{\partial v} \cdot \frac{\partial \beta}{\partial \mu} = \frac{2(\mu-1)}{v^5} - \frac{2\mu}{v^5} = -\frac{2}{v^5}$$

• For us the relevant factor is the magnitude of the Jacobian, which is dominated by the denominator; so we go with $\frac{1}{v^5}$

• We defined v as $(\alpha + \beta)^{\frac{1}{2}}$

\Rightarrow So finally for our prior we have:

$$P(\alpha, \beta) \sim (\alpha + \beta)^{-\frac{5}{2}}$$

\Rightarrow We have everything together for our posterior now!

\Rightarrow We can sample from this using M-H.

\Rightarrow We are in 72 dimensions, it is going to be painful.

\Rightarrow If we can get the conditionals we can do Gibbs sampling!

Reminder: Our posterior is

$$P(\theta, \alpha, \beta | Y, n) \propto P(\alpha, \beta) \cdot \prod_i P(\theta_i | \alpha, \beta) \cdot \prod_i P(Y_i | \theta_i, \alpha, \beta, n_i)$$

$P(\theta_i | Y_i, n_i, \alpha, \beta)$ is easy!

$\Rightarrow P(\alpha, \beta)$ is just a number for given α and β ,
and we have our conjugate prior setup
with a Binomial and a Beta, so we get
another Beta.

$$\Rightarrow P(\theta_i | Y_i, n_i, \alpha, \beta) \propto \text{Beta}(\alpha + Y_i, \beta + n_i - Y_i)$$

We still need the conditionals to sample α and β .

It is not easy, but doable:

$$\frac{\theta_i^{\alpha-1} \cdot (1-\theta_i)^{\beta-1}}{\text{Beta}(\alpha, \beta)} \cdot \underbrace{\theta_i^{Y_i} \cdot (1-\theta_i)^{n_i-Y_i}}_{\text{Binomial (dropped normalization)}}$$

On Wikipedia we read:

$$\Gamma(\alpha) + \Gamma(\beta) = \Gamma(\alpha + \beta) \cdot \text{Beta}(\alpha, \beta)$$

$$\Leftrightarrow \text{B}(\alpha, \beta) = \frac{\Gamma(\alpha) \cdot \Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

This gives us:

$$\frac{\Gamma(\alpha + \beta) \cdot \Theta_i^{\alpha - y_i} \cdot (1 - \Theta_i)^{\beta - 1 + n_i - y_i}}{\Gamma(\alpha) \cdot \Gamma(\beta)}$$

Conditioning on all except α we get:

$$P(\alpha | Y, n, \Theta, \beta) \propto P(\alpha, \beta) \cdot \left(\frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)} \right)^N \cdot \prod_i \Theta_i^\alpha$$

And for β :

$$P(\beta | Y, n, \Theta, \alpha) \propto P(\alpha, \beta) \cdot \left(\frac{\Gamma(\alpha + \beta)}{\Gamma(\beta)} \right)^N \cdot \prod_i (1 - \Theta_i)^\beta$$

Sampling from these conditionals is not straight forward.

\Rightarrow We can still do Gibbs sampling, using M-H to sample from the conditionals for α and β .

\Rightarrow We have a nice sampling function for our Θ_i 's.

The 71st data point:

- In the beginning we had a data point in addition to the data we got from our literature review:

$$n_{71} = 14 \quad Y_{71} = 4$$

- Let's compare how our additional data influences our estimate for θ_{71} compared to the naive estimate $\frac{4}{14}$.

$$P(\theta_{71}, \alpha, \beta | Y_1, n, Y_{71}, n_{71})$$

$$\propto P(\alpha, \beta) \cdot \underbrace{P(\theta_{71} | \alpha, \beta) \cdot P(Y_{71} | \theta_{71}, \alpha, \beta, n_{71})}_{\text{Beta}(\alpha+4, \beta+10)}$$

→ We have samples for this from our Gibbs sampling!