

CPSC 535

Gibbs Sampling

AD

February 2007

- Rejection Sampling and Importance Sampling are two general methods but limited to problems of moderate dimensions.

- Rejection Sampling and Importance Sampling are two general methods but limited to problems of moderate dimensions.
- **Problem:** We try to sample all the components of a potentially high-dimensional parameter simultaneously/sequentially and we can never correct for components already sampled.

- Rejection Sampling and Importance Sampling are two general methods but limited to problems of moderate dimensions.
- **Problem:** We try to sample all the components of a potentially high-dimensional parameter simultaneously/sequentially and we can never correct for components already sampled.
- A powerful class of methods is available to deal with such methods: Markov chain Monte Carlo.

- Multiple failures in a nuclear plant

Pump i	1	2	3	4	5
# Failures p_i	5	1	5	14	3
Times t_i	94.32	15.72	62.88	125.76	5.24
Pump i	6	7	8	9	10
# Failures p_i	19	1	1	4	22
Times t_i	31.44	1.05	1.05	2.10	10.48

- Multiple failures in a nuclear plant

Pump i	1	2	3	4	5
# Failures p_i	5	1	5	14	3
Times t_i	94.32	15.72	62.88	125.76	5.24
Pump i	6	7	8	9	10
# Failures p_i	19	1	1	4	22
Times t_i	31.44	1.05	1.05	2.10	10.48

- Model: Failures of the i -th pump follow a Poisson process with parameter λ_i ($1 \leq i \leq 10$). For an observed time t_i , the number of failures p_i is thus a Poisson $\mathcal{P}(\lambda_i t_i)$ random variable.

- Multiple failures in a nuclear plant

Pump i	1	2	3	4	5
# Failures p_i	5	1	5	14	3
Times t_i	94.32	15.72	62.88	125.76	5.24
Pump i	6	7	8	9	10
# Failures p_i	19	1	1	4	22
Times t_i	31.44	1.05	1.05	2.10	10.48

- Model: Failures of the i -th pump follow a Poisson process with parameter λ_i ($1 \leq i \leq 10$). For an observed time t_i , the number of failures p_i is thus a Poisson $\mathcal{P}(\lambda_i t_i)$ random variable.
- The unknown parameters consist of $\theta = (\lambda_1, \dots, \lambda_{10}, \beta)$.

- Hierarchical model

$$\lambda_i | (\alpha, \beta) \stackrel{\text{iid}}{\sim} \mathcal{G}a(\alpha, \beta) \text{ and } \beta \sim \mathcal{G}a(\gamma, \delta)$$

with $\alpha = 1.8$ and $\gamma = 0.01$ and $\delta = 1$.

- Hierarchical model

$$\lambda_i | (\alpha, \beta) \stackrel{\text{iid}}{\sim} \mathcal{G}a(\alpha, \beta) \text{ and } \beta \sim \mathcal{G}a(\gamma, \delta)$$

with $\alpha = 1.8$ and $\gamma = 0.01$ and $\delta = 1$.

- The posterior distribution is proportional to

$$\begin{aligned} & p(\lambda_{1:10}, \beta | p_{1:10}, t_{1:10}) \\ \propto & \prod_{i=1}^{10} \{ (\lambda_i t_i)^{p_i} \exp(-\lambda_i t_i) \lambda_i^{\alpha-1} \exp(-\beta \lambda_i) \} \beta^{10\alpha} \beta^{\gamma-1} \exp(-\delta \beta) \\ \propto & \prod_{i=1}^{10} \{ \lambda_i^{p_i + \alpha - 1} \exp(-(t_i + \beta) \lambda_i) \} \beta^{10\alpha + \gamma - 1} \exp(-\delta \beta). \end{aligned}$$

- Hierarchical model

$$\lambda_i | (\alpha, \beta) \stackrel{\text{iid}}{\sim} \mathcal{G}a(\alpha, \beta) \text{ and } \beta \sim \mathcal{G}a(\gamma, \delta)$$

with $\alpha = 1.8$ and $\gamma = 0.01$ and $\delta = 1$.

- The posterior distribution is proportional to

$$\begin{aligned} & p(\lambda_{1:10}, \beta | p_{1:10}, t_{1:10}) \\ \propto & \prod_{i=1}^{10} \{ (\lambda_i t_i)^{p_i} \exp(-\lambda_i t_i) \lambda_i^{\alpha-1} \exp(-\beta \lambda_i) \} \beta^{10\alpha} \beta^{\gamma-1} \exp(-\delta\beta) \\ \propto & \prod_{i=1}^{10} \{ \lambda_i^{p_i+\alpha-1} \exp(-(t_i + \beta)\lambda_i) \} \beta^{10\alpha+\gamma-1} \exp(-\delta\beta). \end{aligned}$$

- This multidimensional distribution is rather complex. It is not obvious how the rejection method or importance sampling could be used in this context.

- The conditionals have a familiar form

$$p(\lambda_{1:10} | p_{1:10}, t_{1:10}, \beta) = \prod_{i=1}^{10} p(\lambda_i | p_i, t_i, \beta)$$

where

$$\lambda_i | (\beta, t_i, p_i) \sim \mathcal{G}a(p_i + \alpha, t_i + \beta) \text{ for } 1 \leq i \leq 10,$$

and

$$\beta | (\lambda_1, \dots, \lambda_{10}) \sim \mathcal{G}a(\gamma + 10\alpha, \delta + \sum_{i=1}^{10} \lambda_i).$$

- The conditionals have a familiar form

$$p(\lambda_{1:10} | p_{1:10}, t_{1:10}, \beta) = \prod_{i=1}^{10} p(\lambda_i | p_i, t_i, \beta)$$

where

$$\lambda_i | (\beta, t_i, p_i) \sim \mathcal{G}a(p_i + \alpha, t_i + \beta) \text{ for } 1 \leq i \leq 10,$$

and

$$\beta | (\lambda_1, \dots, \lambda_{10}) \sim \mathcal{G}a(\gamma + 10\alpha, \delta + \sum_{i=1}^{10} \lambda_i).$$

- Instead of directly sampling the vector $\theta = (\lambda_1, \dots, \lambda_{10}, \beta)$ at once, one could suggest sampling it iteratively, starting for example with the λ_i 's for a given guess of β , followed by an update of β given the new samples $\lambda_1, \dots, \lambda_{10}$.

My first Gibbs sampler

- Given a sample, at iteration t , $\theta^t := (\lambda_1^t, \dots, \lambda_{10}^t, \beta^t)$ one could proceed as follows at iteration $t + 1$,
- ① $\lambda_i^{t+1} | (\beta^t, t_i, p_i) \sim \mathcal{G}a(p_i + \alpha, t_i + \beta^t)$ for $1 \leq i \leq 10$,
- ② $\beta^{t+1} | (\lambda_1^{t+1}, \dots, \lambda_{10}^{t+1}) \sim \mathcal{G}a(\gamma + 10\alpha, \delta + \sum_{i=1}^{10} \lambda_i^{t+1})$.

My first Gibbs sampler

- Given a sample, at iteration t , $\theta^t := (\lambda_1^t, \dots, \lambda_{10}^t, \beta^t)$ one could proceed as follows at iteration $t + 1$,
- ① $\lambda_i^{t+1} | (\beta^t, t_i, p_i) \sim \mathcal{G}a(p_i + \alpha, t_i + \beta^t)$ for $1 \leq i \leq 10$,
- ② $\beta^{t+1} | (\lambda_1^{t+1}, \dots, \lambda_{10}^{t+1}) \sim \mathcal{G}a(\gamma + 10\alpha, \delta + \sum_{i=1}^{10} \lambda_i^{t+1})$.
- Instead of directly sampling in a space with 11 dimensions, one samples in spaces of dimension 1.

My first Gibbs sampler

- Given a sample, at iteration t , $\theta^t := (\lambda_1^t, \dots, \lambda_{10}^t, \beta^t)$ one could proceed as follows at iteration $t + 1$,
 - 1 $\lambda_i^{t+1} | (\beta^t, t_i, p_i) \sim \mathcal{G}a(p_i + \alpha, t_i + \beta^t)$ for $1 \leq i \leq 10$,
 - 2 $\beta^{t+1} | (\lambda_1^{t+1}, \dots, \lambda_{10}^{t+1}) \sim \mathcal{G}a(\gamma + 10\alpha, \delta + \sum_{i=1}^{10} \lambda_i^{t+1})$.
- Instead of directly sampling in a space with 11 dimensions, one samples in spaces of dimension 1.
- Note that the deterministic version of such an algorithm where sampling is replaced by maximization would not generally converge towards the global maximum of the joint distribution.

- The structure of the algorithm calls for many questions:

- The structure of the algorithm calls for many questions:
 - Are we sampling from the desired joint distribution?

- The structure of the algorithm calls for many questions:
 - Are we sampling from the desired joint distribution?
 - If yes, how many times should the iteration above be repeated?

- The structure of the algorithm calls for many questions:
 - Are we sampling from the desired joint distribution?
 - If yes, how many times should the iteration above be repeated?
- The validity of the approach described here stems from the fact that the sequence $\{\theta^t\}$ defined above is a Markov chain and some Markov chains have very nice properties.

Elements of Markov chains

- **Markov chain:** A sequence of random variables $\{X_n; n \in \mathbb{N}\}$ defined on $(\mathbb{X}, \mathcal{B}(\mathbb{X}))$ which satisfies the property, for any $A \in \mathcal{B}(\mathbb{X})$

$$\mathbb{P}(X_n \in A | X_0, \dots, X_{n-1}) = \mathbb{P}(X_n \in A | X_{n-1}).$$

and we will write

$$\mathbb{P}(X_n \in A | X_{n-1}) = P(x, A) = \int_A P(x, dy).$$

- **Markov chain:** A sequence of random variables $\{X_n; n \in \mathbb{N}\}$ defined on $(\mathbb{X}, \mathcal{B}(\mathbb{X}))$ which satisfies the property, for any $A \in \mathcal{B}(\mathbb{X})$

$$\mathbb{P}(X_n \in A | X_0, \dots, X_{n-1}) = \mathbb{P}(X_n \in A | X_{n-1}).$$

and we will write

$$\mathbb{P}(X_n \in A | X_{n-1}) = P(x, A) = \int_A P(x, dy).$$

- **Markov chain Monte Carlo:** Given a target π , design a transition kernel P such that asymptotically as $n \rightarrow \infty$

$$\frac{1}{N} \sum_{n=1}^N \varphi(X_n) \rightarrow \int_{\mathbb{X}} \varphi(x) \pi(x) dx \text{ and/or } X_n \sim \pi.$$

Elements of Markov chains

- **Markov chain:** A sequence of random variables $\{X_n; n \in \mathbb{N}\}$ defined on $(\mathbb{X}, \mathcal{B}(\mathbb{X}))$ which satisfies the property, for any $A \in \mathcal{B}(\mathbb{X})$

$$\mathbb{P}(X_n \in A | X_0, \dots, X_{n-1}) = \mathbb{P}(X_n \in A | X_{n-1}).$$

and we will write

$$\mathbb{P}(X_n \in A | X_{n-1}) = P(x, A) = \int_A P(x, dy).$$

- **Markov chain Monte Carlo:** Given a target π , design a transition kernel P such that asymptotically as $n \rightarrow \infty$

$$\frac{1}{N} \sum_{n=1}^N \varphi(X_n) \rightarrow \int_{\mathbb{X}} \varphi(x) \pi(x) dx \text{ and/or } X_n \sim \pi.$$

- It should be easy to simulate the Markov chain even if π is complex.

- Consider the autoregression for $|\alpha| < 1$

$$X_n = \alpha X_{n-1} + V_n, \text{ where } V_n \sim \mathcal{N}(0, \sigma^2)$$

then

$$P(x, dy) = P(x, y) dy = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y - \alpha x)^2}{2\sigma^2}\right) dx.$$

- Consider the autoregression for $|\alpha| < 1$

$$X_n = \alpha X_{n-1} + V_n, \text{ where } V_n \sim \mathcal{N}(0, \sigma^2)$$

then

$$P(x, dy) = P(x, y) dy = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y - \alpha x)^2}{2\sigma^2}\right) dx.$$

- The limiting distribution is

$$\pi(x) = \mathcal{N}\left(x; 0, \frac{\sigma^2}{1 - \alpha^2}\right)$$

and satisfies

$$\int_{\mathbb{X}} \pi(x) P(x, y) dx = \pi(y)$$

- Consider the autoregression for $|\alpha| < 1$

$$X_n = \alpha X_{n-1} + V_n, \text{ where } V_n \sim \mathcal{N}(0, \sigma^2)$$

then

$$P(x, dy) = P(x, y) dy = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y - \alpha x)^2}{2\sigma^2}\right) dx.$$

- The limiting distribution is

$$\pi(x) = \mathcal{N}\left(x; 0, \frac{\sigma^2}{1 - \alpha^2}\right)$$

and satisfies

$$\int_{\mathbb{X}} \pi(x) P(x, y) dx = \pi(y)$$

- To sample from π , we could just sample the Markov chain and asymptotically we would have $X_n \sim \pi$. [Obviously, in this case this is useless because we can sample from π directly.]

- Graphically, consider 1000 independent Markov chains run in parallel.

- Graphically, consider 1000 independent Markov chains run in parallel.
- We assume that the initial distribution of these Markov chains is $\mathcal{U}_{[0,20]}$. So initially, the Markov chains samples are not distributed according to π

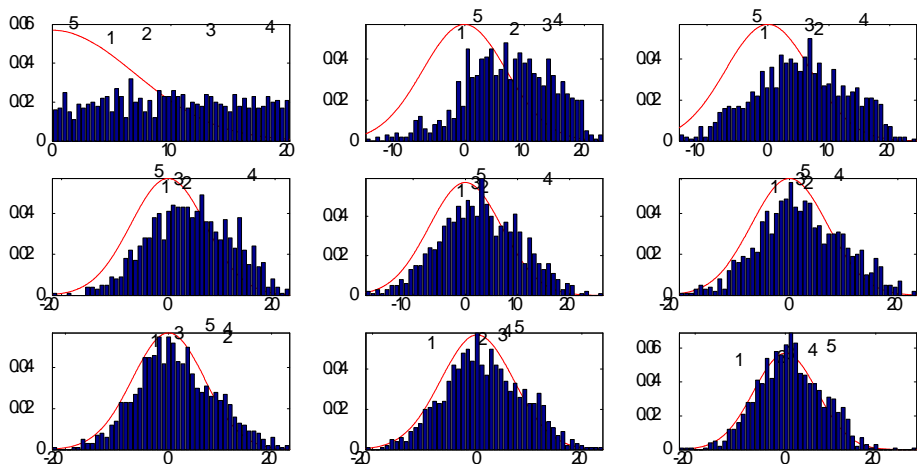


Figure: From top left to bottom right: histograms of 1000 independent Markov chains with a normal distribution as target distribution as n increases.

- The target normal distribution seems to “attract” the distribution of the samples and even to be a fixed point of the algorithm.

- The target normal distribution seems to “attract” the distribution of the samples and even to be a fixed point of the algorithm.
- This is what we wanted to achieve, *i.e.* it seems that we have produced 1000 independent samples from the normal distribution.

- The target normal distribution seems to “attract” the distribution of the samples and even to be a fixed point of the algorithm.
- This is what we wanted to achieve, *i.e.* it seems that we have produced 1000 independent samples from the normal distribution.
- In fact one can show that in many (all?) situations of interest it is not necessary to run N Markov chains in parallel in order to obtain 1000 samples, but that one can consider a unique Markov chain, and build the histogram from this single Markov chain by forming histograms from one trajectory.

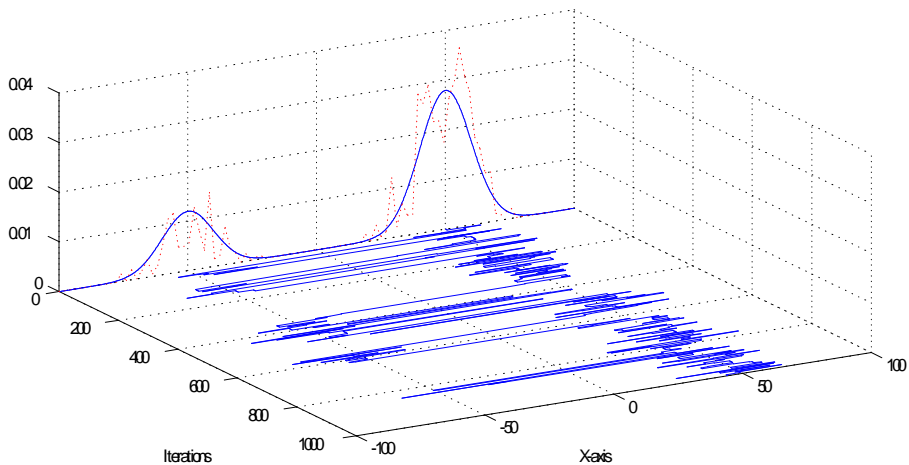


Figure: Bimodal target distributions and simulated Markov chain

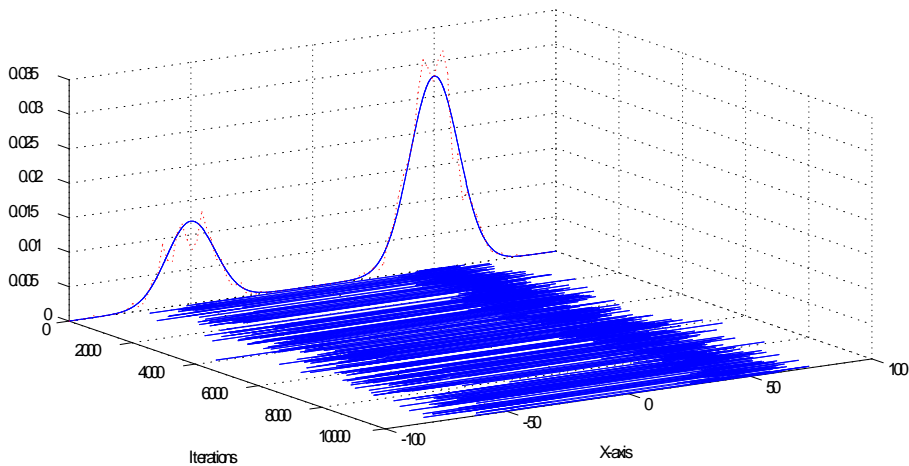


Figure: Bimodal target distributions and simulated Markov chain

- The estimate of the target distribution, through the series of histograms, improves with the number of iterations.

- The estimate of the target distribution, through the series of histograms, improves with the number of iterations.
- Assume that we have stored $\{X_n, 1 \leq n \leq N\}$ for N large and wish to estimate $\int_{\mathbb{X}} \varphi(x) \pi(x) dx$.

- The estimate of the target distribution, through the series of histograms, improves with the number of iterations.
- Assume that we have stored $\{X_n, 1 \leq n \leq N\}$ for N large and wish to estimate $\int_{\mathbb{X}} \varphi(x) \pi(x) dx$.
- In the light of the numerical experiments, one can suggest the estimator

$$\frac{1}{N} \sum_{n=1}^N \varphi(X_n).$$

which is exactly the estimator that we would use if $\{X_n, 1 \leq n \leq N\}$ were independent.

- The estimate of the target distribution, through the series of histograms, improves with the number of iterations.
- Assume that we have stored $\{X_n, 1 \leq n \leq N\}$ for N large and wish to estimate $\int_{\mathbb{X}} \varphi(x) \pi(x) dx$.
- In the light of the numerical experiments, one can suggest the estimator

$$\frac{1}{N} \sum_{n=1}^N \varphi(X_n).$$

which is exactly the estimator that we would use if $\{X_n, 1 \leq n \leq N\}$ were independent.

- In fact, it can be proved, under relatively mild conditions, that such an estimator is consistent *despite the fact that the samples are NOT independent!* Under additional conditions, a CLT also holds with a rate of CV in $1/\sqrt{N}$.

To summarize, we are interested in Markov chains with transition kernel P which have the following three important properties observed above:

- The desired distribution π is a “fixed point” of the algorithm or, in more appropriate terms, an *invariant distribution* of the Markov chain, i.e. $\int_{\mathbb{X}} \pi(x)P(x, y)dx = \pi(y)$.

Markov chains for Monte Carlo

To summarize, we are interested in Markov chains with transition kernel P which have the following three important properties observed above:

- The desired distribution π is a “fixed point” of the algorithm or, in more appropriate terms, an *invariant distribution* of the Markov chain, *i.e.* $\int_{\mathbb{X}} \pi(x)P(x, y)dx = \pi(y)$.
- The successive distributions of the Markov chains converge towards π .

Markov chains for Monte Carlo

To summarize, we are interested in Markov chains with transition kernel P which have the following three important properties observed above:

- The desired distribution π is a “fixed point” of the algorithm or, in more appropriate terms, an *invariant distribution* of the Markov chain, i.e. $\int_{\mathbb{X}} \pi(x)P(x, y)dx = \pi(y)$.
- The successive distributions of the Markov chains converge towards π .
- The estimator $\frac{1}{N} \sum_{n=1}^N \varphi(X_n)$ converges towards $\mathbb{E}_{\pi}(\varphi(X))$ and asymptotically $X_n \sim \pi$

- Given $\pi(x)$, there is an infinite number of kernels $P(x, y)$ which admits $\pi(x)$ as their invariant distribution.

- Given $\pi(x)$, there is an infinite number of kernels $P(x, y)$ which admits $\pi(x)$ as their invariant distribution.
- The “art” of MCMC consists of coming up with good ones.

- Given $\pi(x)$, there is an infinite number of kernels $P(x, y)$ which admits $\pi(x)$ as their invariant distribution.
- The “art” of MCMC consists of coming up with good ones.
- Convergence is ensured under very weak assumptions; namely irreducibility and aperiodicity.

- Given $\pi(x)$, there is an infinite number of kernels $P(x, y)$ which admits $\pi(x)$ as their invariant distribution.
- The “art” of MCMC consists of coming up with good ones.
- Convergence is ensured under very weak assumptions; namely irreducibility and aperiodicity.
- It is usually very easy to establish that an MCMC sampler converges towards π but very difficult to obtain rates of convergence.

Two component Gibbs sampler

- Consider the target distribution $\pi(\theta)$ such that $\theta = (\theta^1, \theta^2)$. Then the 2 component Gibbs sampler proceeds as follows.

Two component Gibbs sampler

- Consider the target distribution $\pi(\theta)$ such that $\theta = (\theta^1, \theta^2)$. Then the 2 component Gibbs sampler proceeds as follows.
- Initialization: Select deterministically or randomly $\theta_0 = (\theta_0^1, \theta_0^2)$.

Two component Gibbs sampler

- Consider the target distribution $\pi(\theta)$ such that $\theta = (\theta^1, \theta^2)$. Then the 2 component Gibbs sampler proceeds as follows.
- Initialization: Select deterministically or randomly $\theta_0 = (\theta_0^1, \theta_0^2)$.
- Iteration i ; $i \geq 1$

Two component Gibbs sampler

- Consider the target distribution $\pi(\theta)$ such that $\theta = (\theta^1, \theta^2)$. Then the 2 component Gibbs sampler proceeds as follows.
- Initialization: Select deterministically or randomly $\theta_0 = (\theta_0^1, \theta_0^2)$.
- Iteration $i; i \geq 1$
 - Sample $\theta_i^1 \sim \pi(\theta^1 | \theta_{i-1}^2)$.

Two component Gibbs sampler

- Consider the target distribution $\pi(\theta)$ such that $\theta = (\theta^1, \theta^2)$. Then the 2 component Gibbs sampler proceeds as follows.
- Initialization: Select deterministically or randomly $\theta_0 = (\theta_0^1, \theta_0^2)$.
- Iteration i ; $i \geq 1$
 - Sample $\theta_i^1 \sim \pi(\theta^1 \mid \theta_{i-1}^2)$.
 - Sample $\theta_i^2 \sim \pi(\theta^2 \mid \theta_i^1)$.

Two component Gibbs sampler

- Consider the target distribution $\pi(\theta)$ such that $\theta = (\theta^1, \theta^2)$. Then the 2 component Gibbs sampler proceeds as follows.
- Initialization: Select deterministically or randomly $\theta_0 = (\theta_0^1, \theta_0^2)$.
- Iteration i ; $i \geq 1$
 - Sample $\theta_i^1 \sim \pi(\theta^1 | \theta_{i-1}^2)$.
 - Sample $\theta_i^2 \sim \pi(\theta^2 | \theta_i^1)$.
- Sampling from these conditional is often feasible even when sampling from the joint is impossible (e.g. nuclear pump data).

- Clearly $\{(\theta_i^1, \theta_i^2)\}$ is a Markov chain and its transition kernel is

$$P\left((\theta^1, \theta^2), (\tilde{\theta}^1, \tilde{\theta}^2)\right) = \pi\left(\tilde{\theta}^1 \mid \theta^2\right) \pi\left(\tilde{\theta}^2 \mid \tilde{\theta}^1\right).$$

- Clearly $\{(\theta_i^1, \theta_i^2)\}$ is a Markov chain and its transition kernel is

$$P\left((\theta^1, \theta^2), (\tilde{\theta}^1, \tilde{\theta}^2)\right) = \pi\left(\tilde{\theta}^1 \mid \theta^2\right) \pi\left(\tilde{\theta}^2 \mid \tilde{\theta}^1\right).$$

- Then $\int \int \pi(\theta^1, \theta^2) P\left((\theta^1, \theta^2), (\tilde{\theta}^1, \tilde{\theta}^2)\right) d\theta^1 d\theta^2$ satisfies

$$\begin{aligned} & \int \int \pi(\theta^1, \theta^2) \pi\left(\tilde{\theta}^1 \mid \theta^2\right) \pi\left(\tilde{\theta}^2 \mid \tilde{\theta}^1\right) d\theta^1 d\theta^2 \\ &= \int \pi(\theta^2) \pi\left(\tilde{\theta}^1 \mid \theta^2\right) \pi\left(\tilde{\theta}^2 \mid \tilde{\theta}^1\right) d\theta^2 \\ &= \int \pi\left(\tilde{\theta}^1, \theta^2\right) \pi\left(\tilde{\theta}^2 \mid \tilde{\theta}^1\right) d\theta^2 \\ &= \pi\left(\tilde{\theta}^1\right) \pi\left(\tilde{\theta}^2 \mid \tilde{\theta}^1\right) = \pi\left(\tilde{\theta}^1, \tilde{\theta}^2\right) \end{aligned}$$

- This does not ensure that the Gibbs sampler does converge towards the invariant distribution!

- This does not ensure that the Gibbs sampler does converge towards the invariant distribution!
- Additionally it is required to ensure *irreducibility*: loosely speaking the Markov chain can move to any set A such that $\pi(A) > 0$ for (almost) any starting point.

- This does not ensure that the Gibbs sampler does converge towards the invariant distribution!
- Additionally it is required to ensure *irreducibility*: loosely speaking the Markov chain can move to any set A such that $\pi(A) > 0$ for (almost) any starting point.
- This ensures that

$$\frac{1}{N} \sum_{n=1}^N \varphi(\theta_n^1, \theta_n^2) \rightarrow \int \varphi(\theta^1, \theta^2) \pi(\theta^1, \theta^2) d\theta^1 d\theta^2$$

but NOT that asymptotically $(\theta_n^1, \theta_n^2) \sim \pi$.

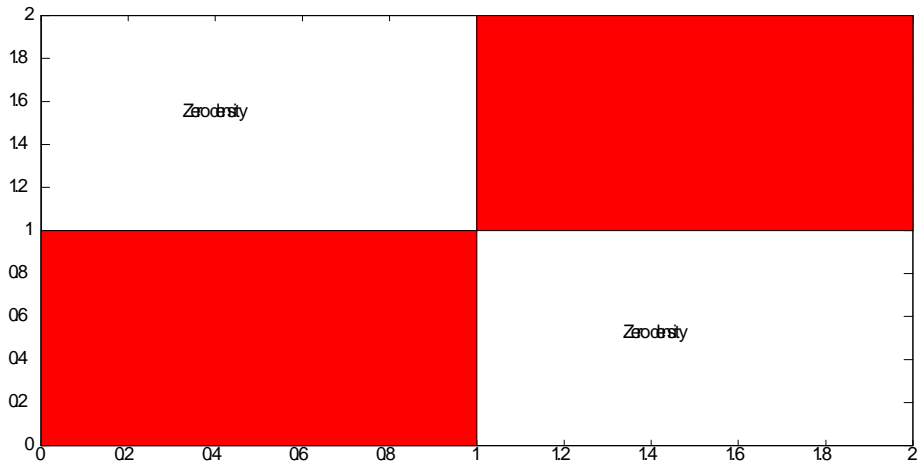


Figure: A distribution that can lead to a reducible Gibbs sampler.

- Consider a simple example where $\mathbb{X} = \{1, 2\}$ and $P(1, 2) = P(2, 1) = 1$. Clearly the invariant distribution is given by $\pi(1) = \pi(2) = \frac{1}{2}$.

- Consider a simple example where $\mathbb{X} = \{1, 2\}$ and $P(1, 2) = P(2, 1) = 1$. Clearly the invariant distribution is given by $\pi(1) = \pi(2) = \frac{1}{2}$.
- However, we know that if the chain starts in $X_0 = 1$, then $X_{2n} = 1$ and $X_{2n+1} = 0$ for any n .

- Consider a simple example where $\mathbb{X} = \{1, 2\}$ and $P(1, 2) = P(2, 1) = 1$. Clearly the invariant distribution is given by $\pi(1) = \pi(2) = \frac{1}{2}$.
- However, we know that if the chain starts in $X_0 = 1$, then $X_{2n} = 1$ and $X_{2n+1} = 0$ for any n .
- We have

$$\frac{1}{N} \sum_{n=1}^N \varphi(X_n) \rightarrow \int \varphi(x) \pi(x) dx$$

but clearly X_n is NOT distributed according to π .

- Consider a simple example where $\mathbb{X} = \{1, 2\}$ and $P(1, 2) = P(2, 1) = 1$. Clearly the invariant distribution is given by $\pi(1) = \pi(2) = \frac{1}{2}$.
- However, we know that if the chain starts in $X_0 = 1$, then $X_{2n} = 1$ and $X_{2n+1} = 0$ for any n .
- We have

$$\frac{1}{N} \sum_{n=1}^N \varphi(X_n) \rightarrow \int \varphi(x) \pi(x) dx$$

but clearly X_n is NOT distributed according to π .

- You need to make sure that you do NOT explore the space in a periodic way to ensure that $X_n \sim \pi$ asymptotically.

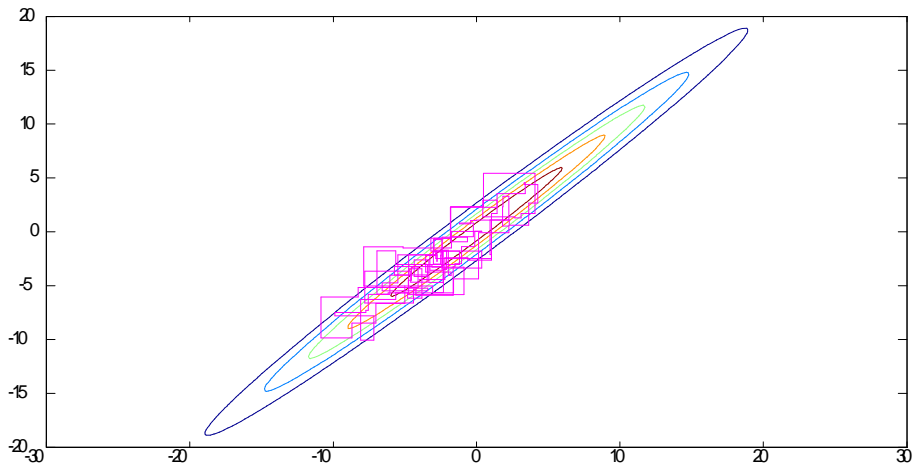


Figure: Even when irreducibility and aperiodicity are ensured, the Gibbs sampler can still converge very slowly.

Deterministic Scan Gibbs Sampler

- If $\theta = (\theta_1, \dots, \theta_p)$ where $p > 2$, the Gibbs sampling strategy still applies.

Deterministic Scan Gibbs Sampler

- If $\theta = (\theta_1, \dots, \theta_p)$ where $p > 2$, the Gibbs sampling strategy still applies.
- Initialization: Select deterministically or randomly $\theta_0 = (\theta_0^1, \dots, \theta_0^p)$.

Deterministic Scan Gibbs Sampler

- If $\theta = (\theta_1, \dots, \theta_p)$ where $p > 2$, the Gibbs sampling strategy still applies.
- Initialization: Select deterministically or randomly $\theta_0 = (\theta_0^1, \dots, \theta_0^p)$.
- Iteration i ; $i \geq 1$:

Deterministic Scan Gibbs Sampler

- If $\theta = (\theta_1, \dots, \theta_p)$ where $p > 2$, the Gibbs sampling strategy still applies.
- Initialization: Select deterministically or randomly $\theta_0 = (\theta_0^1, \dots, \theta_0^p)$.
- Iteration i ; $i \geq 1$:
 - For $k = 1 : p$

Deterministic Scan Gibbs Sampler

- If $\theta = (\theta_1, \dots, \theta_p)$ where $p > 2$, the Gibbs sampling strategy still applies.
- Initialization: Select deterministically or randomly $\theta_0 = (\theta_0^1, \dots, \theta_0^p)$.
- Iteration i ; $i \geq 1$:
 - For $k = 1 : p$
 - Sample $\theta_i^k \sim \pi(\theta^k | \theta_i^{-k})$ where $\theta_i^{-k} = (\theta_i^1, \dots, \theta_i^{k-1}, \theta_{i-1}^{k+1}, \dots, \theta_{i-1}^p)$.

Random Scan Gibbs Sampler

- Initialization: Select deterministically or randomly $\theta_0 = (\theta_0^1, \dots, \theta_0^p)$.

Random Scan Gibbs Sampler

- Initialization: Select deterministically or randomly $\theta_0 = (\theta_0^1, \dots, \theta_0^p)$.
- Iteration i ; $i \geq 1$:

Random Scan Gibbs Sampler

- Initialization: Select deterministically or randomly $\theta_0 = (\theta_0^1, \dots, \theta_0^p)$.
- Iteration i ; $i \geq 1$:
 - Sample $K \sim U_{\{1, \dots, p\}}$.

Random Scan Gibbs Sampler

- Initialization: Select deterministically or randomly $\theta_0 = (\theta_0^1, \dots, \theta_0^p)$.
- Iteration i ; $i \geq 1$:
 - Sample $K \sim U_{\{1, \dots, p\}}$.
 - Set $\theta_i^{-K} = \theta_{i-1}^{-K}$.

Random Scan Gibbs Sampler

- Initialization: Select deterministically or randomly $\theta_0 = (\theta_0^1, \dots, \theta_0^p)$.
- Iteration i ; $i \geq 1$:
 - Sample $K \sim U_{\{1, \dots, p\}}$.
 - Set $\theta_i^{-K} = \theta_{i-1}^{-K}$.
 - Sample $\theta_i^K \sim \pi(\theta^K \mid \theta_i^{-K})$ where
 $\theta_i^{-K} = (\theta_i^1, \dots, \theta_i^{K-1}, \theta_i^{K+1}, \dots, \theta_i^p)$.

Tricks of the trade

- Try to have as few “blocks” as possible.

Tricks of the trade

- Try to have as few “blocks” as possible.
- Put the most correlated variables in the same block.

Tricks of the trade

- Try to have as few “blocks” as possible.
- Put the most correlated variables in the same block.
- If necessary, reparametrize the model to achieve this.

Tricks of the trade

- Try to have as few “blocks” as possible.
- Put the most correlated variables in the same block.
- If necessary, reparametrize the model to achieve this.
- Integrate analytically as many variables as possible: pretty algorithms can be much more inefficient than ugly algorithms.

Tricks of the trade

- Try to have as few “blocks” as possible.
- Put the most correlated variables in the same block.
- If necessary, reparametrize the model to achieve this.
- Integrate analytically as many variables as possible: pretty algorithms can be much more inefficient than ugly algorithms.
- There is no general result telling strategy A is better than strategy B in all cases: you need experience.

Application to Simulation of Fractal Images

- Consider a 2D black and white 'target' image. We define an distribution ν which assigns $1/P$ mass on each black point and zero on white points where P is the number of black points.

Application to Simulation of Fractal Images

- Consider a 2D black and white 'target' image. We define an distribution ν which assigns $1/P$ mass on each black point and zero on white points where P is the number of black points.
- Now we consider the following simple Markov process on \mathbb{R}^2 with

$$P(x, y) = \sum_{i=1}^k w_i \delta_{A_i x + b_i}(y)$$

and we select $\{w_i, A_i, b_i\}$ such that $P(x, dy)$ has an invariant distribution π which is an approximation of ν .

Application to Simulation of Fractal Images

- Consider a 2D black and white 'target' image. We define an distribution ν which assigns $1/P$ mass on each black point and zero on white points where P is the number of black points.
- Now we consider the following simple Markov process on \mathbb{R}^2 with

$$P(x, y) = \sum_{i=1}^k w_i \delta_{A_i x + b_i}(y)$$

and we select $\{w_i, A_i, b_i\}$ such that $P(x, dy)$ has an invariant distribution π which is an approximation of ν .

- To find $\{w_i, A_i, b_i\}$, we write

$$\begin{aligned} \int \pi(x) P(x, y) f(y) dx dy &= \sum_{i=1}^k w_i \int f(A_i x + b_i) \pi(x) dx \\ &= \int f(x) \pi(x) dx \approx \int f(x) \nu(x) dx \end{aligned}$$

and solve approximately the equations for some functions f (linear or low order polynomials).

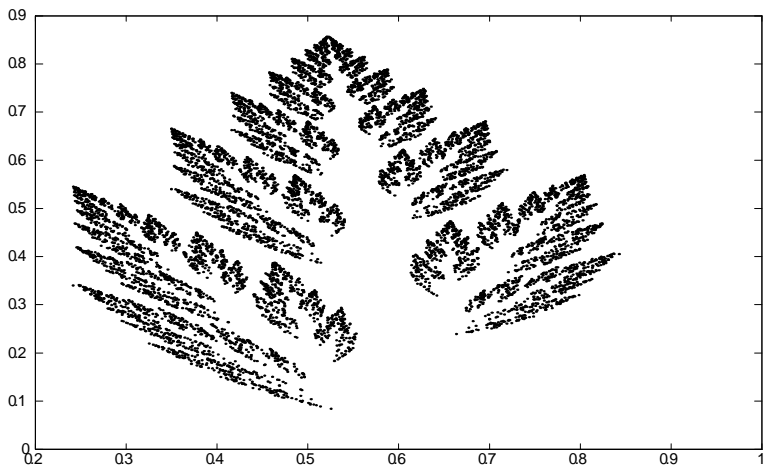


Figure: Fractal image generated using iterated random functions with $k = 2$ and $N = 10000$ samples

Gibbs Sampler for Bayesian Variable Selection

- We select the following model

$$Y = \sum_{i=1}^p \beta_i X_i + \sigma V \text{ where } V \sim \mathcal{N}(0, 1)$$

where we assume $\mathcal{IG}(\sigma^2; \frac{\nu_0}{2}, \frac{\gamma_0}{2})$ and for $\alpha^2 \ll 1$

$$\beta_i \sim \frac{1}{2} \mathcal{N}(0, \alpha^2 \delta^2 \sigma^2) + \frac{1}{2} \mathcal{N}(0, \delta^2 \sigma^2)$$

Gibbs Sampler for Bayesian Variable Selection

- We select the following model

$$Y = \sum_{i=1}^p \beta_i X_i + \sigma V \text{ where } V \sim \mathcal{N}(0, 1)$$

where we assume $\mathcal{IG}(\sigma^2; \frac{\nu_0}{2}, \frac{\gamma_0}{2})$ and for $\alpha^2 \ll 1$

$$\beta_i \sim \frac{1}{2} \mathcal{N}(0, \alpha^2 \delta^2 \sigma^2) + \frac{1}{2} \mathcal{N}(0, \delta^2 \sigma^2)$$

- We introduce a latent variable $\gamma_i \in \{0, 1\}$ such that

$$\Pr(\gamma_i = 0) = \Pr(\gamma_i = 1) = \frac{1}{2},$$
$$\beta_i | \gamma_i = 0 \sim \mathcal{N}(0, \alpha^2 \delta^2 \sigma^2), \quad \beta_i | \gamma_i = 1 \sim \mathcal{N}(0, \delta^2 \sigma^2).$$

- We have parameters $(\beta_{1:p}, \gamma_{1:p}, \sigma^2)$ and observe $D = \{x_i, y_i\}_{i=1}^n$.

- We have parameters $(\beta_{1:p}, \gamma_{1:p}, \sigma^2)$ and observe $D = \{x_i, y_i\}_{i=1}^n$.
- A potential Gibbs sampler consists of sampling iteratively from $p(\beta_{1:p} | D, \gamma_{1:p}, \sigma^2)$ (Gaussian), $p(\sigma^2 | D, \gamma_{1:p}, \beta_{1:p})$ (inverse-Gamma) and $p(\gamma_{1:p} | D, \beta_{1:p}, \sigma^2)$.

- We have parameters $(\beta_{1:p}, \gamma_{1:p}, \sigma^2)$ and observe $D = \{x_i, y_i\}_{i=1}^n$.
- A potential Gibbs sampler consists of sampling iteratively from $p(\beta_{1:p} | D, \gamma_{1:p}, \sigma^2)$ (Gaussian), $p(\sigma^2 | D, \gamma_{1:p}, \beta_{1:p})$ (inverse-Gamma) and $p(\gamma_{1:p} | D, \beta_{1:p}, \sigma^2)$.
- In particular

$$p(\gamma_{1:p} | D, \beta_{1:p}, \sigma^2) = \prod_{i=1}^p p(\gamma_i | \beta_i, \sigma^2)$$

and

$$p(\gamma_i = 1 | \beta_i, \sigma^2) = \frac{\frac{1}{\sqrt{2\pi\delta\sigma}} \exp\left(-\frac{\beta_i^2}{2\delta^2\sigma^2}\right)}{\frac{1}{\sqrt{2\pi\delta\sigma}} \exp\left(-\frac{\beta_i^2}{2\delta^2\sigma^2}\right) + \frac{1}{\sqrt{2\pi\alpha\delta\sigma}} \exp\left(-\frac{\beta_i^2}{2\alpha^2\delta^2\sigma^2}\right)}.$$

- We have parameters $(\beta_{1:p}, \gamma_{1:p}, \sigma^2)$ and observe $D = \{x_i, y_i\}_{i=1}^n$.
- A potential Gibbs sampler consists of sampling iteratively from $p(\beta_{1:p} | D, \gamma_{1:p}, \sigma^2)$ (Gaussian), $p(\sigma^2 | D, \gamma_{1:p}, \beta_{1:p})$ (inverse-Gamma) and $p(\gamma_{1:p} | D, \beta_{1:p}, \sigma^2)$.
- In particular

$$p(\gamma_{1:p} | D, \beta_{1:p}, \sigma^2) = \prod_{i=1}^p p(\gamma_i | \beta_i, \sigma^2)$$

and

$$p(\gamma_i = 1 | \beta_i, \sigma^2) = \frac{\frac{1}{\sqrt{2\pi\delta\sigma}} \exp\left(-\frac{\beta_i^2}{2\delta^2\sigma^2}\right)}{\frac{1}{\sqrt{2\pi\delta\sigma}} \exp\left(-\frac{\beta_i^2}{2\delta^2\sigma^2}\right) + \frac{1}{\sqrt{2\pi\alpha\delta\sigma}} \exp\left(-\frac{\beta_i^2}{2\alpha^2\delta^2\sigma^2}\right)}.$$

- The Gibbs sampler becomes reducible as α goes to zero.

- This is the result of bad modelling and bad algorithm. You would like to put $\alpha \simeq 0$ and write

$$Y = \sum_{i=1}^p \gamma_i \beta_i X_i + \sigma V \text{ where } V \sim \mathcal{N}(0, 1)$$

where $\gamma_i = 1$ if X_i is included or $\gamma_i = 0$ otherwise. However this suggests that β_i is defined even when $\gamma_i = 0$.

- This is the result of bad modelling and bad algorithm. You would like to put $\alpha \simeq 0$ and write

$$Y = \sum_{i=1}^p \gamma_i \beta_i X_i + \sigma V \text{ where } V \sim \mathcal{N}(0, 1)$$

where $\gamma_i = 1$ if X_i is included or $\gamma_i = 0$ otherwise. However this suggests that β_i is defined even when $\gamma_i = 0$.

- A neater way to write such models is to write

$$Y = \sum_{\{i:\gamma_i=1\}} \beta_i X_i + \sigma V = \beta_\gamma^\top X_\gamma + \sigma V$$

where, for a vector $\gamma = (\gamma_1, \dots, \gamma_p)$, $\beta_\gamma = \{\beta_i : \gamma_i = 1\}$, $X_\gamma = \{X_i : \gamma_i = 1\}$ and $n_\gamma = \sum_{i=1}^p \gamma_i$.

- This is the result of bad modelling and bad algorithm. You would like to put $\alpha \simeq 0$ and write

$$Y = \sum_{i=1}^p \gamma_i \beta_i X_i + \sigma V \text{ where } V \sim \mathcal{N}(0, 1)$$

where $\gamma_i = 1$ if X_i is included or $\gamma_i = 0$ otherwise. However this suggests that β_i is defined even when $\gamma_i = 0$.

- A neater way to write such models is to write

$$Y = \sum_{\{i:\gamma_i=1\}} \beta_i X_i + \sigma V = \beta_\gamma^\top X_\gamma + \sigma V$$

where, for a vector $\gamma = (\gamma_1, \dots, \gamma_p)$, $\beta_\gamma = \{\beta_i : \gamma_i = 1\}$, $X_\gamma = \{X_i : \gamma_i = 1\}$ and $n_\gamma = \sum_{i=1}^p \gamma_i$.

- Prior distributions

$$\pi_\gamma(\beta_\gamma, \sigma^2) = \mathcal{N}(\beta_\gamma; 0, \delta^2 \sigma^2 I_{n_\gamma}) \mathcal{IG}(\sigma^2; \frac{\nu_0}{2}, \frac{\gamma_0}{2})$$

and $\pi(\gamma) = \prod_{i=1}^p \pi(\gamma_i) = 2^{-p}$.

- We are interested in sampling from the trans-dimensional distribution $\pi(\gamma, \beta_\gamma, \sigma^2 \mid D)$.

- We are interested in sampling from the trans-dimensional distribution $\pi(\gamma, \beta_\gamma, \sigma^2 | D)$.
- However, we know that

$$\pi(\gamma, \beta_\gamma, \sigma^2 | D) = \pi(\gamma | D) \pi(\beta_\gamma, \sigma^2 | D, \gamma)$$

where

$$\pi(\gamma | D) \propto \pi(D | \gamma) \pi(\gamma)$$

and

$$\pi(D | \gamma) = \int \pi(D, \beta_\gamma, \sigma^2 | \gamma) d\beta_\gamma d\sigma^2.$$

- $\pi(\gamma|D)$ is a discrete probability distribution with 2^p potential values. We can use the Gibbs sampler to sample from it.

- $\pi(\gamma|D)$ is a discrete probability distribution with 2^p potential values. We can use the Gibbs sampler to sample from it.
- Initialization: Select deterministically or randomly $\gamma_0 = (\gamma_0^1, \dots, \gamma_0^p)$.

- $\pi(\gamma|D)$ is a discrete probability distribution with 2^p potential values. We can use the Gibbs sampler to sample from it.
- Initialization: Select deterministically or randomly $\gamma_0 = (\gamma_0^1, \dots, \gamma_0^p)$.
- Iteration i ; $i \geq 1$:

- $\pi(\gamma | D)$ is a discrete probability distribution with 2^p potential values. We can use the Gibbs sampler to sample from it.
- Initialization: Select deterministically or randomly $\gamma_0 = (\gamma_0^1, \dots, \gamma_0^p)$.
- Iteration $i; i \geq 1$:
 - For $k = 1 : p$

- $\pi(\gamma | D)$ is a discrete probability distribution with 2^p potential values. We can use the Gibbs sampler to sample from it.
- Initialization: Select deterministically or randomly $\gamma_0 = (\gamma_0^1, \dots, \gamma_0^p)$.
- Iteration $i; i \geq 1$:
 - For $k = 1 : p$
 - Sample $\gamma_i^k \sim \pi(\gamma_i^k | D, \gamma_i^{-k})$ where $\gamma_i^{-k} = (\gamma_i^1, \dots, \gamma_i^{k-1}, \gamma_i^{k+1}, \dots, \gamma_i^p)$.

- $\pi(\gamma | D)$ is a discrete probability distribution with 2^p potential values. We can use the Gibbs sampler to sample from it.
- Initialization: Select deterministically or randomly $\gamma_0 = (\gamma_0^1, \dots, \gamma_0^p)$.
- Iteration $i; i \geq 1$:
 - For $k = 1 : p$
 - Sample $\gamma_i^k \sim \pi(\gamma^k | D, \gamma_i^{-k})$ where $\gamma_i^{-k} = (\gamma_i^1, \dots, \gamma_i^{k-1}, \gamma_{i-1}^{k+1}, \dots, \gamma_{i-1}^p)$.
 - Optional step: Sample $(\beta_{\gamma,i}, \sigma_i^2) \sim \pi(\beta_{\gamma}, \sigma^2 | D, \gamma)$.

- This very simple sampler is much more efficient than the previous one.

- This very simple sampler is much more efficient than the previous one.
- However, it can also mix very slowly because the components are updated one at a time.

- This very simple sampler is much more efficient than the previous one.
- However, it can also mix very slowly because the components are updated one at a time.
- Updating correlated components together would increase significantly the convergence speed of the algorithm at the cost of an increased complexity.

Finite Mixture Models

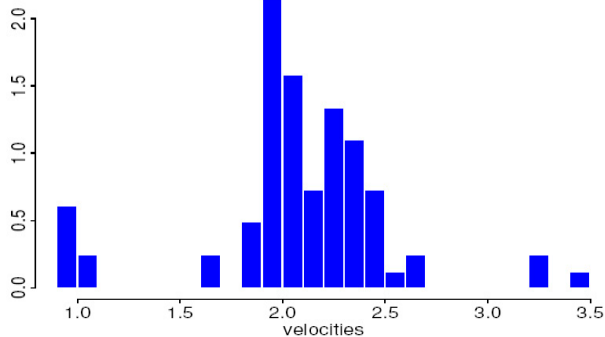


Figure: Velocity (km/sc) of galaxies in the Corona Borealis Region

- Consider the case where one has n data X_i

$$X_i \stackrel{\text{i.i.d}}{\sim} \sum_{k=1}^K p_k \mathcal{N}(\mu_k, \sigma_k^2)$$

where K is fixed and $\theta = \{\mu_k, \sigma_k^2, p_k\}_{k=1, \dots, K}$ are unknown.

- Consider the case where one has n data X_i

$$X_i \stackrel{\text{i.i.d.}}{\sim} \sum_{k=1}^K p_k \mathcal{N}(\mu_k, \sigma_k^2)$$

where K is fixed and $\theta = \{\mu_k, \sigma_k^2, p_k\}_{k=1, \dots, K}$ are unknown.

- A standard approach consists of finding a local maximum of the log-likelihood

$$\log f(x_{1:n} | \theta) = \sum_{i=1}^n \log f(x_i | \theta)$$

where

$$f(x | \theta) = \sum_{k=1}^K \frac{p_k}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{(x - \mu_k)^2}{2\sigma_k^2}\right).$$

- Consider the case where one has n data X_i

$$X_i \stackrel{\text{i.i.d.}}{\sim} \sum_{k=1}^K p_k \mathcal{N}(\mu_k, \sigma_k^2)$$

where K is fixed and $\theta = \{\mu_k, \sigma_k^2, p_k\}_{k=1, \dots, K}$ are unknown.

- A standard approach consists of finding a local maximum of the log-likelihood

$$\log f(x_{1:n} | \theta) = \sum_{i=1}^n \log f(x_i | \theta)$$

where

$$f(x | \theta) = \sum_{k=1}^K \frac{p_k}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{(x - \mu_k)^2}{2\sigma_k^2}\right).$$

- *Problem:* The likelihood is unbounded!

- We consider the Bayesian framework where we set priors

$$\pi(\theta) = \pi(p_1, \dots, p_K) \prod_{k=1}^K \pi(\mu_k, \sigma_k^2).$$

- We consider the Bayesian framework where we set priors

$$\pi(\theta) = \pi(p_1, \dots, p_K) \prod_{k=1}^K \pi(\mu_k, \sigma_k^2).$$

- We use the following (conditionally conjugate) priors where

$$(p_1, \dots, p_K) \sim \mathcal{D}(\gamma_1, \dots, \gamma_K).$$
$$\mu_k | \sigma_k^2 \sim \mathcal{N}\left(\alpha_k, \frac{\sigma_k^2}{\lambda_k}\right), \quad \sigma_k^2 \sim \mathcal{IG}\left(\frac{\lambda_k + 3}{2}, \frac{\beta_k}{2}\right).$$

- We consider the Bayesian framework where we set priors

$$\pi(\theta) = \pi(p_1, \dots, p_K) \prod_{k=1}^K \pi(\mu_k, \sigma_k^2).$$

- We use the following (conditionally conjugate) priors where

$$(p_1, \dots, p_K) \sim \mathcal{D}(\gamma_1, \dots, \gamma_K).$$
$$\mu_k | \sigma_k^2 \sim \mathcal{N}\left(\alpha_k, \frac{\sigma_k^2}{\lambda_k}\right), \sigma_k^2 \sim \mathcal{IG}\left(\frac{\lambda_k + 3}{2}, \frac{\beta_k}{2}\right).$$

- It is impossible to use the Gibbs sampler to sample from $\pi(\theta | x_{1:n})$.

- Like in the EM, we can introduce the missing data $Z_i \in \{1, \dots, K\}$ such that

$$X_i | Z_i \sim \mathcal{N}(\mu_{Z_i}, \sigma_{Z_i}^2)$$

and

$$\Pr(Z_i = k) = p_k.$$

- Like in the EM, we can introduce the missing data $Z_i \in \{1, \dots, K\}$ such that

$$X_i | Z_i \sim \mathcal{N}(\mu_{Z_i}, \sigma_{Z_i}^2)$$

and

$$\Pr(Z_i = k) = p_k.$$

- The “complete” likelihood admits a simple form

$$\pi(x_{1:n}, z_{1:n} | \theta) = \prod_{k=1}^n f(x_i | \theta, z_i) \pi(z_i | \theta).$$

- Like in the EM, we can introduce the missing data $Z_i \in \{1, \dots, K\}$ such that

$$X_i | Z_i \sim \mathcal{N}(\mu_{Z_i}, \sigma_{Z_i}^2)$$

and

$$\Pr(Z_i = k) = p_k.$$

- The “complete” likelihood admits a simple form

$$\pi(x_{1:n}, z_{1:n} | \theta) = \prod_{k=1}^n f(x_i | \theta, z_i) \pi(z_i | \theta).$$

- Thus we propose to sample the joint posterior $\pi(\theta, z_{1:n} | y_{1:n})$ using the Gibbs sampler; that is sampling iteratively from $\pi(\theta | y_{1:n}, z_{1:n})$ and $\pi(z_{1:n} | y_{1:n}, \theta)$.

Gibbs Sampler for Finite Mixture Distributions

- We have

$$\pi(z_{1:n} | \theta, x_{1:n}) = \prod_{i=1}^n \pi(z_i | \theta, x_i)$$

where

$$\pi(z_i = j | \theta, x_i) = \frac{f(x_i | \theta, j) p_j}{\sum_{k=1}^K f(x_i | \theta, k) p_k}.$$

Gibbs Sampler for Finite Mixture Distributions

- We have

$$\pi(z_{1:n} | \theta, x_{1:n}) = \prod_{i=1}^n \pi(z_i | \theta, x_i)$$

where

$$\pi(z_i = j | \theta, x_i) = \frac{f(x_i | \theta, j) p_j}{\sum_{k=1}^K f(x_i | \theta, k) p_k}.$$

- We have

$$\pi(\theta | z_{1:n}, x_{1:n}) = \pi(p_1, \dots, p_K | z_{1:n}) \prod_{k=1}^K \pi(\mu_k, \sigma_k^2 | z_{1:n}, x_{1:n})$$

- Introducing

$$n_k = \sum_{i=1}^n \mathbf{1}_{\{k\}}(z_i), n_k \bar{x}_k = \sum_{i=1}^n x_i \mathbf{1}_{\{k\}}(z_i), s_k^2 = \sum_{i=1}^n (x_i - \bar{x}_k)^2 \mathbf{1}_{\{k\}}(z_i).$$

- Introducing

$$n_k = \sum_{i=1}^n \mathbf{1}_{\{k\}}(z_i), n_k \bar{x}_k = \sum_{i=1}^n x_i \mathbf{1}_{\{k\}}(z_i), s_k^2 = \sum_{i=1}^n (x_i - \bar{x}_k)^2 \mathbf{1}_{\{k\}}(z_i).$$

- We have the full conditionals

$$p_1, \dots, p_K | z_{1:n} \sim \mathcal{D}(\gamma_1 + n_1, \dots, \gamma_K + n_K),$$

$$\sigma_k^2 | z_{1:n}, x_{1:n} \sim \text{IG} \left(\frac{\lambda_k + n_k + 3}{2}, \frac{\lambda_k s_k^2 + \beta_k + s_k^2 - (\lambda_k + n_k)^{-1} (\lambda_k \alpha_k + n_k \bar{x}_k)^2}{2} \right),$$

$$\mu_k | \sigma_k^2, z_{1:n}, x_{1:n} \sim \mathcal{N} \left(\frac{\lambda_k \alpha_k + n_k \bar{x}_k}{\lambda_k + n_k}, \frac{\sigma_k^2}{\lambda_k + n_k} \right).$$

- Introducing

$$n_k = \sum_{i=1}^n \mathbf{1}_{\{k\}}(z_i), n_k \bar{x}_k = \sum_{i=1}^n x_i \mathbf{1}_{\{k\}}(z_i), s_k^2 = \sum_{i=1}^n (x_i - \bar{x}_k)^2 \mathbf{1}_{\{k\}}(z_i).$$

- We have the full conditionals

$$p_1, \dots, p_K | z_{1:n} \sim \mathcal{D}(\gamma_1 + n_1, \dots, \gamma_K + n_K),$$

$$\sigma_k^2 | z_{1:n}, x_{1:n} \sim \text{IG} \left(\frac{\lambda_k + n_k + 3}{2}, \frac{\lambda_k s_k^2 + \beta_k + s_k^2 - (\lambda_k + n_k)^{-1} (\lambda_k \alpha_k + n_k \bar{x}_k)^2}{2} \right),$$

$$\mu_k | \sigma_k^2, z_{1:n}, x_{1:n} \sim \mathcal{N} \left(\frac{\lambda_k \alpha_k + n_k \bar{x}_k}{\lambda_k + n_k}, \frac{\sigma_k^2}{\lambda_k + n_k} \right).$$

- It is thus trivial to implement the Gibbs sampler.

Simulation Results

- Consider some $n = 100$ simulated data

$$X_i \sim 0.3\mathcal{N}(-2, 1) + 0.7\mathcal{N}(2, 1),$$

i.e. we have well-separated components.

Simulation Results

- Consider some $n = 100$ simulated data

$$X_i \sim 0.3\mathcal{N}(-2, 1) + 0.7\mathcal{N}(2, 1),$$

i.e. we have well-separated components.

- We set $\gamma_k = 1$, $\alpha_k = 0$, $\lambda_k = 0.01$, $\beta_k = 0.01$ and run the Gibbs sampler for 10000 iterations.

- Consider some $n = 100$ simulated data

$$X_i \sim 0.3\mathcal{N}(-2, 1) + 0.7\mathcal{N}(2, 1),$$

i.e. we have well-separated components.

- We set $\gamma_k = 1$, $\alpha_k = 0$, $\lambda_k = 0.01$, $\beta_k = 0.01$ and run the Gibbs sampler for 10000 iterations.
- We obtain $\hat{\mathbb{E}}(\mu_1 | x_{1:n}) = 2.17$, $\hat{\mathbb{E}}(\mu_2 | x_{1:n}) = -1.89$,
 $\hat{\mathbb{E}}(\sigma_1^2 | x_{1:n}) = 0.92$, $\hat{\mathbb{E}}(\sigma_2^2 | x_{1:n}) = 1.3$, $\hat{\mathbb{E}}(p_1 | x_{1:n}) = 0.32$ and
 $\hat{\mathbb{E}}(p_2 | x_{1:n}) = 0.68$.

Simulation Results

- Consider some $n = 100$ simulated data

$$X_i \sim 0.3\mathcal{N}(-2, 1) + 0.7\mathcal{N}(2, 1),$$

i.e. we have well-separated components.

- We set $\gamma_k = 1$, $\alpha_k = 0$, $\lambda_k = 0.01$, $\beta_k = 0.01$ and run the Gibbs sampler for 10000 iterations.
- We obtain $\hat{\mathbb{E}}(\mu_1 | x_{1:n}) = 2.17$, $\hat{\mathbb{E}}(\mu_2 | x_{1:n}) = -1.89$,
 $\hat{\mathbb{E}}(\sigma_1^2 | x_{1:n}) = 0.92$, $\hat{\mathbb{E}}(\sigma_2^2 | x_{1:n}) = 1.3$, $\hat{\mathbb{E}}(p_1 | x_{1:n}) = 0.32$ and
 $\hat{\mathbb{E}}(p_2 | x_{1:n}) = 0.68$.
- Increasing the number of iterations to 100000, I obtain similar results.
Any good?

- Your algorithm does not work! Indeed we know that

$$\begin{aligned}\mathbb{E}(\mu_1 | x_{1:n}) &= \mathbb{E}(\mu_2 | x_{1:n}), \quad \mathbb{E}(\sigma_1^2 | x_{1:n}) = \mathbb{E}(\sigma_2^2 | x_{1:n}), \\ \mathbb{E}(p_1 | x_{1:n}) &= \mathbb{E}(p_2 | x_{1:n}) = 0.5.\end{aligned}$$

- Your algorithm does not work! Indeed we know that

$$\begin{aligned}\mathbb{E}(\mu_1 | x_{1:n}) &= \mathbb{E}(\mu_2 | x_{1:n}), \quad \mathbb{E}(\sigma_1^2 | x_{1:n}) = \mathbb{E}(\sigma_2^2 | x_{1:n}), \\ \mathbb{E}(p_1 | x_{1:n}) &= \mathbb{E}(p_2 | x_{1:n}) = 0.5.\end{aligned}$$

- This follows because both the prior and likelihood are exchangeable, that is

$$\begin{aligned}&\pi(p_1, \dots, p_K, \mu_1, \dots, \mu_K, \sigma_1^2, \dots, \sigma_K^2 | x_{1:n}) \\ &= \pi(p_{\zeta(1)}, \dots, p_{\zeta(K)}, \mu_{\zeta(1)}, \dots, \mu_{\zeta(K)}, \sigma_{\zeta(1)}^2, \dots, \sigma_{\zeta(K)}^2 | x_{1:n})\end{aligned}$$

for any permutation ζ of the labels.

- Your algorithm does not work! Indeed we know that

$$\begin{aligned}\mathbb{E}(\mu_1 | x_{1:n}) &= \mathbb{E}(\mu_2 | x_{1:n}), \quad \mathbb{E}(\sigma_1^2 | x_{1:n}) = \mathbb{E}(\sigma_2^2 | x_{1:n}), \\ \mathbb{E}(p_1 | x_{1:n}) &= \mathbb{E}(p_2 | x_{1:n}) = 0.5.\end{aligned}$$

- This follows because both the prior and likelihood are exchangeable, that is

$$\begin{aligned}&\pi(p_1, \dots, p_K, \mu_1, \dots, \mu_K, \sigma_1^2, \dots, \sigma_K^2 | x_{1:n}) \\ &= \pi(p_{\zeta(1)}, \dots, p_{\zeta(K)}, \mu_{\zeta(1)}, \dots, \mu_{\zeta(K)}, \sigma_{\zeta(1)}^2, \dots, \sigma_{\zeta(K)}^2 | x_{1:n})\end{aligned}$$

for any permutation ζ of the labels.

- Clearly, conditional expectations are not useful in this case.
 \Rightarrow This does NOT mean that your Bayesian model is useless.

- One can select another point estimates; e.g. the MAP estimate

$$\theta_{MAP} = \arg \max \pi(\theta | x_{1:n}).$$

- One can select another point estimates; e.g. the MAP estimate

$$\theta_{MAP} = \arg \max \pi(\theta | x_{1:n}).$$

- Alternatively, constraints can be set on the priors; e.g. we ensure that

$$\mu_1 \leq \mu_2 \leq \dots \leq \mu_P$$

- One can select another point estimates; e.g. the MAP estimate

$$\theta_{MAP} = \arg \max \pi(\theta | x_{1:n}).$$

- Alternatively, constraints can be set on the priors; e.g. we ensure that

$$\mu_1 \leq \mu_2 \leq \dots \leq \mu_P$$

- However, this can lead to “strange” shapes of the posteriors and is not natural in most cases.

- One way to improve the algorithm consists of randomly permuting the labels (Fruwirth-Schnatter, JASA, 2002)

⇒ Realistic only if K is moderate because there are $K!$ permutations.

- One way to improve the algorithm consists of randomly permuting the labels (Fruwirth-Schnatter, JASA, 2002)

⇒ Realistic only if K is moderate because there are $K!$ permutations.

- Alternative ways to improve the algorithm include

- One way to improve the algorithm consists of randomly permuting the labels (Fruwirth-Schnatter, JASA, 2002)

⇒ Realistic only if K is moderate because there are $K!$ permutations.

- Alternative ways to improve the algorithm include
 - Not introducing the latent variables and using sampling strategies different from Gibbs.

- One way to improve the algorithm consists of randomly permuting the labels (Fruwirth-Schnatter, JASA, 2002)

⇒ Realistic only if K is moderate because there are $K!$ permutations.

- Alternative ways to improve the algorithm include
 - Not introducing the latent variables and using sampling strategies different from Gibbs.
 - Integrating out θ as the marginal distribution $\pi(z_{1:n} | x_{1:n})$ can be computed analytically (for conjugate priors)

- Initialization: Select deterministically or randomly $z_{1:n}^{(0)}$.

- Initialization: Select deterministically or randomly $z_{1:n}^{(0)}$.
- Iteration i ; $i \geq 1$:

- Initialization: Select deterministically or randomly $z_{1:n}^{(0)}$.
- Iteration i ; $i \geq 1$:
 - For $k = 1 : n$, sample $Z_k^{(i)} \sim \pi \left(z_k \mid x_{1:n}, z_{-k}^{(i)} \right)$ where

$$z_{-k}^{(i)} = \left(z_1^{(i)}, \dots, z_{k-1}^{(i)}, z_{k+1}^{(i-1)}, \dots, z_n^{(i-1)} \right).$$

- Initialization: Select deterministically or randomly $z_{1:n}^{(0)}$.
- Iteration i ; $i \geq 1$:
 - For $k = 1 : n$, sample $Z_k^{(i)} \sim \pi \left(z_k \mid x_{1:n}, z_{-k}^{(i)} \right)$ where

$$z_{-k}^{(i)} = \left(z_1^{(i)}, \dots, z_{k-1}^{(i)}, z_{k+1}^{(i-1)}, \dots, z_n^{(i-1)} \right).$$
 - Sample $\theta^{(i)} \sim \pi \left(\theta \mid x_{1:n}, z_{1:n}^{(i)} \right)$.

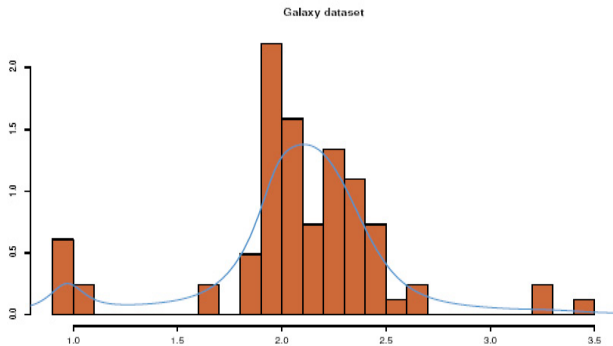


Figure: Predictive distribution for the galaxy dataset.

- The Gibbs sampler is a generic tool to sample approximately from high-dimensional distributions.

- The Gibbs sampler is a generic tool to sample approximately from high-dimensional distributions.
- Each time you face a problem, you need to think hard about it to design an efficient algorithm.

- The Gibbs sampler is a generic tool to sample approximately from high-dimensional distributions.
- Each time you face a problem, you need to think hard about it to design an efficient algorithm.
- Except the choice of the partitions of parameters, the Gibbs sampler is parameter free; this does not mean it is efficient.