

Markov Chain Monte Carlo

Charlie Murry

Boston College

April 11, 2019

Overview of Bayesian Estimation

- We wish to know about unknown parameter $\theta^0 \in R^N$
- We have data y , and $L(y|\theta)$ is the likelihood of y given that $\theta = \theta^0$.

Frequentist

- Derive an estimator (MLE) and analyze statistical properties of that estimator

$$\hat{\theta} = \max_{\theta} L(y|\theta).$$

Bayesian

- Start with a prior belief, $p(\theta)$.
- Use data to update their belief to posterior using Bayes Rule

$$\pi(\theta|y) = \frac{L(y|\theta)p(\theta)}{f(y)}$$

where $f(y) = \int L(y|\theta)\pi(\theta)d\theta$ is the marginal distribution of y .

Bayes Rule

$$\pi(\theta|y) = \frac{L(y|\theta)p(\theta)}{f(y)}$$

is just

$$Pr(A|B) = \frac{Pr(A \cap B)}{Pr(B)} = \frac{Pr(B|A)Pr(A)}{Pr(B)}.$$

Bayesian Estimation Overview

- Outcome of estimation is $\rho(\cdot)$: summarizes everything we know about where θ is.
- Typically report moments of ρ .
- ρ is typically not tractable. How do we report moments from something that is not tractable.

Good: No need to solve complex optimization problem.

Bad: By construction there is a complex integral.

Tractable Example I

Cameron and Trivedi, p422

Nothing about a posterior per se that requires MCMC.

- Suppose you observe N draws from a normal distribution with mean θ and variance σ^2 , e.g.,

$$y_i \sim N(\theta, \sigma^2).$$

- σ^2 is known, but you want to estimate θ .
- A frequentist might use maximum likelihood estimation. The likelihood is:

$$\begin{aligned} L(y|\theta) &= \prod_{i=1}^N (2\pi\sigma^2)^{-\frac{1}{2}} \exp \left\{ -\frac{(y_i - \theta)^2}{2\sigma^2} \right\} \\ &= (2\pi\sigma^2)^{-\frac{N}{2}} \exp \left\{ -\sum_{i=1}^N \frac{(y_i - \theta)^2}{2\sigma^2} \right\} \\ &\propto \exp \left\{ -\frac{N}{2\sigma^2} (\bar{y} - \theta)^2 \right\} \end{aligned}$$

Tractable Example II

Cameron and Trivedi, p422

- Clearly this is maximized at \bar{y} , which is the MLE estimate.
- Just compute the average from your data.

Bayesian

- Define prior belief, θ .
- Suppose that belief is normally distributed with mean μ and variance τ^2
- Prior density:

$$p(\theta) = (2\pi\tau^2)^{-\frac{1}{2}} \exp \left\{ -\frac{(\theta - \mu)^2}{2\tau^2} \right\}.$$

Tractable Example III

Cameron and Trivedi, p422

Following Bayes Rule, the Posterior is proportional to:

$$\begin{aligned}\pi(\theta|y) &\propto L(y|\theta)p(\theta) \\ &\propto \exp\left\{-\frac{N}{2\sigma^2}(\bar{y} - \theta)^2\right\} \exp\left\{-\frac{(\theta - \mu)^2}{2\tau^2}\right\} \\ &\propto \exp\left\{-\frac{1}{2}\left[\frac{(\theta - \tilde{\mu})^2}{\tilde{\tau}^2}\right]\right\}\end{aligned}$$

Where,

$$\begin{aligned}\tilde{\mu} &= \tilde{\tau}^2 \left(\frac{N}{\sigma^2} \bar{y} + \frac{1}{\tau^2} \mu \right) \\ \tilde{\tau}^2 &= \left(\frac{N}{\sigma^2} + \frac{1}{\tau^2} \right)^{-1}\end{aligned}$$

Tractable Example III

Cameron and Trivedi, p422

- The final line is a normal kernel (just complete the square :)).
- The posterior is normally distributed with mean $\tilde{\mu}$ which is a weighted sum of y and the prior mean μ .
- Since this posterior is normal, it is easy for us to compute the moments.
- Mean of posterior goes to \bar{y} as $N \rightarrow \infty$.
- But computing moments of even slightly messier posteriors will require complex integration that we will tackle via simulation.

Review of Monte Carlo Integration

The point of monte carlo integration is to use draws from a distribution to calculate the moments of $\rho(\theta|y)$. If $\rho(\cdot)$ is “easy” to draw from (say, uniform or normal) than we can use traditional monte carlo integration techniques:

$$E[m(\theta)] = \int_{\Theta} m(\theta)\rho(\theta|y)d\theta \approx \frac{1}{S} \sum_{s=1}^S m(\theta_s)$$

- $m(\cdot)$ is an arbitrary function, say identity if we want the mean. We need to assume this expectation exists (of course).
- θ_s is a draw from $\rho(\theta|y)$.

However, this isn't helpful if we don't know how to generate draws from $\rho(\cdot)$ and if we did, we could probably just integrate it directly.

Markov Chain Monte Carlo

MCMC uses draws from a **Markov Chain**, instead of i.i.d. draws from some known distribution.

Use MCMC when:

- Analytic solutions aren't tractable.
- IID sampling doesn't give adequate coverage (perhaps dimension is too high or good approximation of ρ is unknown).

The goal becomes constructing an **ergodic** Markov Chain F (so that the **stationary distribution** exists) such that the stationary distribution is exactly ρ . If we do this then we can generate moments of ρ from

$$E[m(\theta)] \approx \frac{1}{S} \sum_{i=1}^S m(\theta_i)$$

where $\theta_i \sim F(\cdot | \theta_{i-1})$.

English, please?

ergodic: statistical properties can be deduced from a single, sufficiently long, random sample of the process.

note ergodic: a process that changes erratically at an inconsistent rate

stationary distribution: probability distribution that remains unchanged in the Markov chain as time progresses. (The transition matrix of a discrete processes remains constant)

Markov Chain Theory

let the state space for θ be discrete, $\Theta = \{\theta^{(1)}, \dots, \theta^{(K)}\}$.¹

Let our chain be defined by

$$P(\theta_{r+1} = \theta^{(j)} | \theta_r = \theta^{(i)}) = p_{ij}$$

So the Markov transition matrix is,

$$P = \begin{bmatrix} p_{11} & p_{12} & \dots & p_{1K} \\ p_{21} & p_{22} & \dots & p_{2K} \\ \vdots & & & \vdots \\ p_{K1} & p_{K2} & \dots & p_{KK} \end{bmatrix}$$

¹Of course this isn't reasonable for estimation but it allows me to skip over a bunch of measure theory.

MC Theory: Stationarity

Let π_0 be an initial distribution over states (a $1 \times K$ vector). Then the distribution over states after 1 period will be:

$$Pr(\theta_1 = \theta^{(j)}) = \sum_{i=1}^K Pr(\theta_0 = \theta^{(i)})p_{ij} = \sum_{i=1}^K \pi_{0i}p_{ij}$$

Or in matrix notation for the entire distribution,

$$\pi_1 = \pi_0 P$$

If for all i, j : $p_{ij} > 0$, then every state will be visited infinitely often.

A stationary distribution exists and is unique:

$$\lim_{r \rightarrow \infty} \pi_0 P^r = \pi$$

for any π_0 and of course,

$$\pi = \pi P$$

The stationary distribution π is sometimes called the invariant distribution.

Time Reversibility

Definition

A chain is time reversible with respect to π if it has the same behavior backwards and forwards starting from π . That is if the chance of seeing a transition from i to j is the same as seeing a transition from j to i :

$$\pi_i p_{ij} = \pi_j p_{ji}$$

MCMC: Gibbs Sampling

Construct Markov chain by “cycling” through conditional distributions related to π .

- Let $\theta = [\theta_1, \theta_2]'$ with posterior density $p(\theta_1, \theta_2)$.
- *If the conditional densities are known* alternating sequential draws from $p(\theta_1 \mid \theta_2)$ and $p(\theta_2 \mid \theta_1)$ converge to $p(\theta_1, \theta_2)$.

Gibbs Example: Probit

Using Data Augmentation

Model:

$$z_i = x_i \beta + \epsilon_i$$

$$y_i = \begin{cases} 0 & z_i \leq 0 \\ 1 & z_i > 0 \end{cases}$$

$$\epsilon_i \sim N(0, 1)$$

We observe a random sample of (y_i, x_i) and want to estimate β .

Suppose we have a prior $\beta \sim N(\bar{\beta}, A^{-1})$. If we observed z_i then the posterior would be normal (normal is the *conjugate prior* of normal).

However, when z is unobserved there is no simple conjugate prior.

Instead, we can use an “augmentation step” by employing a Gibbs sampler with two blocks (z_i, β) , the second step uses draws of z and the normal conjugate prior.

Probit Example — Algorithm

1. Given β_{r-1} , draw z_i by drawing from a truncated normal:

$$z_{i,r} | \beta_{r-1}, y_i, x_i \sim \text{TruncatedNormal}_a^b(-x_i \beta_{r-1}, 1)$$

Where bounds are $a = 0, b = \infty$ if $y_i = 1$ and $a = -\infty, b = 0$ if $y_i = 0$

2. Draw $\beta_r | z_{i,r}, x_i$ from the posterior of a regression of z on x :

$$\beta_r \sim N(\tilde{\beta}, (X'X + A)^{-1})$$

where $\tilde{\beta} = (X'X + A)^{-1}(X'z + A\tilde{\beta})$.

3. After many draws, we have a sample of β_r which we use as draws from the stationary distribution.

Example — Multivariate Normal

The file `simpleGibbs.m` implements Gibbs sampling to draw from a bivariate normal:

$$(y_1, y_2)' \sim N \left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix} \right)$$

This trivially implies conditional distributions:

$$y_1|y_2 \sim N(\mu_1 + \rho \frac{\sigma_1}{\sigma_2}(y_2 - \mu_2), \sigma_1^2(1 - \rho^2))$$