

# Metropolis and Metropolis-Hastings Algorithms

Jingchen (Monika) Hu

Vassar College

MATH 347 Bayesian Statistics

# Outline

- 1 Overview
- 2 Metropolis Algorithm
- 3 Metropolis-Hastings Algorithm
- 4 Summary

# Outline

- 1 Overview
- 2 Metropolis Algorithm
- 3 Metropolis-Hastings Algorithm
- 4 Summary

# Overview

- Not all parameters have recognizable full conditional posterior distributions.
  - ▶ If you use a non-conjugate prior distribution for a parameter, e.g. Normal for  $\mu$  but Uniform for  $\phi$  in the Normal sampling model.
- What to do when parameters do not have recognizable full conditional posterior distributions?

# Overview

- Not all parameters have recognizable full conditional posterior distributions.
  - ▶ If you use a non-conjugate prior distribution for a parameter, e.g. Normal for  $\mu$  but Uniform for  $\phi$  in the Normal sampling model.
- What to do when parameters do not have recognizable full conditional posterior distributions? JAGS!
- But what does JAGS do?

# Overview

- Not all parameters have recognizable full conditional posterior distributions.
  - ▶ If you use a non-conjugate prior distribution for a parameter, e.g. Normal for  $\mu$  but Uniform for  $\phi$  in the Normal sampling model.
- What to do when parameters do not have recognizable full conditional posterior distributions? JAGS!
- But what does JAGS do?
- Two important MCMC techniques: the Metropolis algorithm and the Metropolis-Hastings algorithm

# Outline

- 1 Overview
- 2 Metropolis Algorithm**
- 3 Metropolis-Hastings Algorithm
- 4 Summary

# Metropolis Algorithm

Suppose we want to estimate  $\pi(\theta | Y)$  for some scalar  $\theta$ .

- ① Start with an initial guess at  $\theta$ , say  $\theta^{(1)}$ .
- ② Given  $\theta^{(s)}$ , generate a value  $\theta^{(s+1)}$  as follows:
  - ▶ Draw plausible value of  $\theta$  from some symmetric distribution  $J(\theta | \theta^{(s)})$  that is easy to simulate, like a  $\text{Normal}(\theta^{(s)}, c)$ , i.e.,

$$\theta^* \sim J(\theta | \theta^{(s)}). \quad (1)$$

- ▶ If  $\theta^*$  is more likely under  $\pi(\theta | Y)$  than  $\theta^{(s)}$ , then we keep it as a plausible value of  $\theta$ , i.e.,  $\theta^{(s+1)} = \theta^*$ .
  - ★ If  $\theta^*$  is less likely under  $\pi(\theta | Y)$  than  $\theta^{(s)}$ , then we let  $\theta^{(s+1)} = \theta^*$  with probability

$$r = \frac{\pi(\theta^* | Y)}{\pi(\theta^{(s)} | Y)} = \frac{p(Y | \theta^*)\pi(\theta^*)}{p(Y | \theta^{(s)})\pi(\theta^{(s)})}. \quad (2)$$

- ③ Repeat Step 2 until MCMC convergence (or for a large number of iterations, say  $S = 10^5$ ).



# Features of Jumping Distribution

- $J(\theta \mid \theta^{(s)})$  is called the proposal distribution.
- $J(\theta \mid \theta^{(s)})$  must depend only on  $\theta^{(s)}$  and not previous values of  $\theta$  in the chain.
- $J(\theta \mid \theta^{(s)})$  must be a symmetric density, i.e.,

$$J(\theta^{(s+1)} \mid \theta^{(s)}) = J(\theta^{(s)} \mid \theta^{(s+1)}). \quad (3)$$

- $J(\theta \mid \theta^{(s)})$  must be such that you can get to any value of the parameter space for  $\theta$  eventually from any  $\theta^{(s)}$ .
- $J(\theta \mid \theta^{(s)})$  must be such that you don't return periodically to any particular value of  $\theta$ .

# Tuning Metropolis Algorithm

- You get to specify  $J(\theta \mid \theta^{(s)})$ , e.g., proposal variance.
- Small proposal steps: high acceptance rate, but the moves are never very large so the Markov chain is sticky and highly correlated.
- Large proposal steps: quickly moves to posterior mode but gets “stuck” for long periods, since proposed values are usually far away from the mode.

# Tuning Metropolis Algorithm cont'd

- Goal is to select one that leads to roughly 35% of new proposed  $\theta^{(s+1)}$  accepted (or at least between 20% to 50%).
- Tuning: try short runs and record percentage of acceptances, and reset  $J$  as necessary to achieve near 35%.
- For example, with a Normal jumping distribution, reset the variance  $c^2$  (or standard deviation  $c$ ) until you get about 35% acceptances.

# Metropolis example: Normal-Normal model with known $\sigma$

- The sampling density:

$$y_1, \dots, y_n \mid \mu, \sigma \stackrel{i.i.d.}{\sim} \text{Normal}(\mu, \sigma). \quad (4)$$

- The prior distribution:

$$\mu \sim \text{Normal}(\mu_0, \sigma_0). \quad (5)$$

- The analytical posterior distribution:

$$\mu \mid y_1, \dots, y_n, \phi \sim \text{Normal} \left( \frac{\phi_0 \mu_0 + n \phi \bar{y}}{\phi_0 + n \phi}, \sqrt{\frac{1}{\phi_0 + n \phi}} \right). \quad (6)$$

With values of  $\bar{y}$ ,  $n$ ,  $\phi$  (i.e.  $\sigma$ ),  $\mu_0$ ,  $\phi_0$  (i.e.  $\sigma_0$ ), we know exactly what this posterior distribution is, and we can use Monte Carlo simulation to generate draws of  $\mu$ .

# Metropolis example: Normal-Normal model with known $\sigma$

- The sampling density:

$$y_1, \dots, y_n \mid \mu, \sigma \stackrel{i.i.d.}{\sim} \text{Normal}(\mu, \sigma). \quad (4)$$

- The prior distribution:

$$\mu \sim \text{Normal}(\mu_0, \sigma_0). \quad (5)$$

- The analytical posterior distribution:

$$\mu \mid y_1, \dots, y_n, \phi \sim \text{Normal} \left( \frac{\phi_0 \mu_0 + n \phi \bar{y}}{\phi_0 + n \phi}, \sqrt{\frac{1}{\phi_0 + n \phi}} \right). \quad (6)$$

With values of  $\bar{y}$ ,  $n$ ,  $\phi$  (i.e.  $\sigma$ ),  $\mu_0$ ,  $\phi_0$  (i.e.  $\sigma_0$ ), we know exactly what this posterior distribution is, and we can use Monte Carlo simulation to generate draws of  $\mu$ .

- How about using the Metropolis algorithm to obtain draws of  $\mu$ ?

# Metropolis example: Normal-Normal model with known $\sigma$ cont'd

Choose a Uniform jumping distribution:

$$\mu^* \sim J(\mu \mid \mu^{(s)}) = \text{Uniform}(\mu^{(s)} - C, \mu^{(s)} + C). \quad (7)$$

- Step 1: choose a new value  $\mu^*$  from  $\text{Uniform}(\mu^{(s)} - C, \mu^{(s)} + C)$ .
- Step 2: calculate the ratio:

$$r = \frac{\pi(\mu^* \mid Y)}{\pi(\mu^{(s)} \mid Y)} = \frac{p(Y \mid \mu^*)\pi(\mu^*)}{p(Y \mid \mu^{(s)})\pi(\mu^{(s)})}. \quad (8)$$

How can one compute  $r$ ?

# Metropolis example: Normal-Normal model with known $\sigma$ cont'd

Choose a Uniform jumping distribution:

$$\mu^* \sim J(\mu \mid \mu^{(s)}) = \text{Uniform}(\mu^{(s)} - C, \mu^{(s)} + C). \quad (7)$$

- Step 1: choose a new value  $\mu^*$  from  $\text{Uniform}(\mu^{(s)} - C, \mu^{(s)} + C)$ .
- Step 2: calculate the ratio:

$$r = \frac{\pi(\mu^* \mid Y)}{\pi(\mu^{(s)} \mid Y)} = \frac{p(Y \mid \mu^*)\pi(\mu^*)}{p(Y \mid \mu^{(s)})\pi(\mu^{(s)})}. \quad (8)$$

How can one compute  $r$ ?

$$r = \left( \frac{\prod_i \text{dnorm}(y_i, \mu^*, \sigma)}{\prod_i \text{dnorm}(y_i, \mu^{(s)}, \sigma)} \right) \left( \frac{\text{dnorm}(\mu^*, \mu_0, \sigma_0)}{\text{dnorm}(\mu^{(s)}, \mu_0, \sigma_0)} \right) \quad (9)$$

# Metropolis example: Normal-Normal model with known $\sigma$ cont'd

$$r = \left( \frac{\prod_i \text{dnorm}(y_i, \mu^*, \sigma)}{\prod_i \text{dnorm}(y_i, \mu^{(s)}, \sigma)} \right) \left( \frac{\text{dnorm}(\mu^*, \mu_0, \sigma_0)}{\text{dnorm}(\mu^{(s)}, \mu_0, \sigma_0)} \right) \quad (10)$$

In many cases, computing the ratio  $r$  directly can be numerically unstable. Therefore, one can work with  $\log(r)$ .

$$\begin{aligned} \log(r) = & \sum_i \left( \log \text{dnorm}(y_i, \mu^*, \sigma) - \log \text{dnorm}(y_i, \mu^{(s)}, \sigma) \right) + \\ & \left( \log \text{dnorm}(\mu^*, \mu_0, \sigma_0) - \log \text{dnorm}(\mu^{(s)}, \mu_0, \sigma_0) \right). \end{aligned} \quad (11)$$



# Metropolis example: Normal-Normal model with known $\sigma$

## cont'd

```

llik = sum(dnorm(y, mu, sigma, log = TRUE));

for (t in (thin + 1):iter){
  mup = runif(1, mu - C, mu + C);
  llikp = sum(dnorm(y, mup, sigma, log = TRUE));

  logr = llikp - llik + dnorm(mup, mu0, sigma0, log = TRUE) -
        dnorm(mu, mu0, sigma0, log = TRUE);

  logu = log(runif(1));
  if(logr > logu){
    mu = mup;
    llik = llikp;
    acc0 = acc0 + 1;
  }
}

```

# Metropolis example: q-Gaussian model

- This is a current independent study, extending a previous MATH 347 project.
- q-Gaussian distribution for the sampling model for  $y_1, \dots, y_n$ :

$$p(y_i \mid \mu_q, \sigma_q) = \frac{1}{\sigma_q B(\frac{\alpha}{2}, \frac{1}{2})} \sqrt{\frac{|Z|}{u^{(1+1/Z)}}} \quad (12)$$

where  $Z = (q - 1)/(3 - q)$  and  $a = 1 - 1/Z$  if  $q < 1$ ;  $a = 1/Z$  if  $1 < q < 3$ , and  $u(y_i) = 1 + Z(y_i - \mu_q)^2/\sigma_q^2$ .

- Prior distributions:

$$q \sim \text{Uniform}(0, 5/3) \quad (13)$$

$$\mu_q \sim \text{Normal}(0, 100) \quad (14)$$

$$\sigma_q \sim \text{Uniform}(0, 100) \quad (15)$$

# Metropolis example: q-Gaussian model cont'd

```

qp = runif(1, q - Cq, q + Cq)
  if (qp <= 1 || qp >= 3) {
    qvals[length(qvals) + 1] = q
  } else {
    llikep = likelihood(x, qp, sigma, mu)

    r = (llikep / llike) * dunif(qp, 1, 5/3) / dunif(q, 1, 5/3)
    if (!is.nan(r)) {

      u = runif(1)
      if (r > u) {
        q = qp
        qvals[length(qvals) + 1] = q
        llike = llikep
        acceptedq = acceptedq + 1
      } else {qvals[length(qvals) + 1] = q}
    } else {qvals[length(qvals) + 1] = q}
  }

```

# Outline

- 1 Overview
- 2 Metropolis Algorithm
- 3 Metropolis-Hastings Algorithm**
- 4 Summary

# Motivation

- Sometimes drawing from symmetric proposal distribution  $J(\theta \mid \theta^{(s)})$  not efficient, i.e., takes long time for chain to converge.
- Example of such inefficiency:
  - ▶ Suppose  $\pi(\theta \mid Y)$  has long tail like a Gamma distribution.
  - ▶ Normal proposal with small variance: takes long time to traverse distribution repeatedly.
  - ▶ Normal proposal with large variance: many proposed  $\theta$  with small posterior density, so too small rate of acceptance.

## Motivation cont'd

- In such cases, ideal to propose values in tail roughly in same proportion as they appear in  $\pi(\theta \mid Y)$ .
- For example,  $J \sim \text{Gamma}$  might be a closer approximation to  $\pi(\theta \mid Y)$  than  $J \sim \text{Normal}$ .
- But, Gamma distribution is not symmetric proposal distribution.
- We have to correct the acceptance ratio  $r$  for this fact; otherwise, we might inaccurately favor values with high density in  $J$  that may not be high density in  $p(\theta \mid Y)$ .
- This leads to the Metropolis-Hastings (M-H) algorithm.

# Metropolis-Hastings Algorithm

Suppose we want to estimate  $\pi(\theta \mid Y)$  using M-H

- Propose a new  $\theta^* \sim J(\theta \mid \theta^{(s)})$  where  $J$  is an arbitrary distribution (certain restrictions apply).

- Compute Metropolis-Hastings ratio

$$\alpha = \min \left\{ 1, \frac{\pi(\theta^* \mid Y) J(\theta^{(s)} \mid \theta^*)}{\pi(\theta^{(s)} \mid Y) J(\theta^* \mid \theta^{(s)})} \right\}$$

- Set

$$\theta^{(s+1)} = \begin{cases} \theta^* & \text{with probability } \alpha \\ \theta^{(s)} & \text{with probability } 1 - \alpha \end{cases}$$

# Metropolis-Hastings Algorithm

Suppose we want to estimate  $\pi(\theta \mid Y)$  using M-H

- Propose a new  $\theta^* \sim J(\theta \mid \theta^{(s)})$  where  $J$  is an arbitrary distribution (certain restrictions apply).

- Compute Metropolis-Hastings ratio

$$\alpha = \min \left\{ 1, \frac{\pi(\theta^* \mid Y) J(\theta^{(s)} \mid \theta^*)}{\pi(\theta^{(s)} \mid Y) J(\theta^* \mid \theta^{(s)})} \right\}$$

- Set

$$\theta^{(s+1)} = \begin{cases} \theta^* & \text{with probability } \alpha \\ \theta^{(s)} & \text{with probability } 1 - \alpha \end{cases}$$

- Recall in Metropolis algorithm, the ratio is  $r = \frac{p(\theta^*|Y)}{p(\theta^{(s)}|Y)}$ .

- Think about  $\frac{J(\theta^{(s)}|\theta^*)}{J(\theta^*|\theta^{(s)})}$  as a correction factor.

If symmetric distribution, then this correction factor turns out to be 1



# Features of M-H Jumping Distribution

- It is easy to sample from  $J(\theta \mid \theta^{(s)})$  and to compute  $\alpha$ .
- $J(\theta \mid \theta^{(s)})$  must depend only on  $\theta^{(s)}$  and not previous values of  $\theta$  in the chain.
- $J(\theta \mid \theta^{(s)})$  must be such that you can get to any value of the parameter space for  $\theta$  eventually from any  $\theta^{(s)}$ .
- $J(\theta \mid \theta^{(s)})$  must be such that you don't return periodically to any particular value of  $\theta$ .
- You get to specify  $J(\theta \mid \theta^{(s)})$ . Use tuning to select one that leads to roughly 35% of new proposed  $\theta^*$  accepted.
- Can use different jumping distributions in different iterations, i.e.,  $J$  is allowed to depend on  $s$ . But  $J$  cannot depend on the draws, i.e.,  $\theta^{(s)}$ .

## Special Cases of M-H Algorithm

- Metropolis algorithm: symmetric jump  $J(\theta^* | \theta^{(s)}) = J(\theta^{(s)} | \theta^*)$

$$\alpha = \min \left\{ 1, \frac{\pi(\theta^* | Y)}{\pi(\theta^{(s)} | Y)} \right\}$$

- Gibbs sampler: jumping distribution equals the target distribution, i.e.,  $J(\theta^* | \theta^{(s)}) = \pi(\theta^* | Y)$ , hence

$$\alpha = \min \left\{ 1, \frac{\pi(\theta^* | Y)p(\theta^{(s)} | Y)}{\pi(\theta^{(s)} | Y)p(\theta^* | Y)} \right\} = 1$$

- Since  $J$  can be different in different iterations, we can update each dimension of the parameter vector one at a time, using either Gibbs, Metropolis, or M-H update.

# M-H example: Normal-Normal model with known $\sigma$

- The sampling density:

$$y_1, \dots, y_n \mid \mu, \sigma \stackrel{i.i.d.}{\sim} \text{Normal}(\mu, \sigma). \quad (16)$$

- The prior distribution:

$$\mu \sim \text{Normal}(\mu_0, \sigma_0). \quad (17)$$

- The analytical posterior distribution:

$$\mu \mid y_1, \dots, y_n, \phi \sim \text{Normal} \left( \frac{\phi_0 \mu_0 + n \phi \bar{y}}{\phi_0 + n \phi}, \sqrt{\frac{1}{\phi_0 + n \phi}} \right). \quad (18)$$

With values of  $\bar{y}$ ,  $n$ ,  $\phi$  (i.e.  $\sigma$ ),  $\mu_0$ ,  $\phi_0$  (i.e.  $\sigma_0$ ), we know exactly what this posterior distribution is, and we can use Monte Carlo simulation to generate draws of  $\mu$ .

- How about using the Metropolis-Hastings algorithm to obtain draws of  $\mu$ ?

# M-H example: Normal-Normal model with known $\sigma$ cont'd

Choose a Gamma jumping distribution:

$$\mu^* \sim J(\mu \mid \mu^{(s)}) = \text{Gamma}(\mu^{(s)}, 1). \quad (19)$$

- Step 1: choose a new value  $\mu^*$  from  $\text{Gamma}(\mu^{(s)}, 1)$ .
- Step 2: calculate the ratio:

$$\alpha = \min \left\{ 1, \frac{\pi(\mu^* \mid Y) J(\mu^{(s)} \mid \mu^*)}{\pi(\mu^{(s)} \mid Y) J(\mu^* \mid \mu^{(s)})} \right\} = \min \left\{ 1, \frac{p(Y \mid \mu^*) \pi(\mu^*) J(\mu^{(s)} \mid \mu^*)}{p(Y \mid \mu^{(s)}) \pi(\mu^{(s)}) J(\mu^* \mid \mu^{(s)})} \right\}. \quad (20)$$

How can one compute the ratio  $\alpha^*$ ?

# M-H example: Normal-Normal model with known $\sigma$ cont'd

Choose a Gamma jumping distribution:

$$\mu^* \sim J(\mu \mid \mu^{(s)}) = \text{Gamma}(\mu^{(s)}, 1). \quad (19)$$

- Step 1: choose a new value  $\mu^*$  from  $\text{Gamma}(\mu^{(s)}, 1)$ .
- Step 2: calculate the ratio:

$$\alpha = \min \left\{ 1, \frac{\pi(\mu^* \mid Y) J(\mu^{(s)} \mid \mu^*)}{\pi(\mu^{(s)} \mid Y) J(\mu^* \mid \mu^{(s)})} \right\} = \min \left\{ 1, \frac{p(Y \mid \mu^*) \pi(\mu^*) J(\mu^{(s)} \mid \mu^*)}{p(Y \mid \mu^{(s)}) \pi(\mu^{(s)}) J(\mu^* \mid \mu^{(s)})} \right\}. \quad (20)$$

How can one compute the ratio  $\alpha^*$ ?

$$\alpha^* = \left( \frac{\prod_i \text{dnorm}(y_i, \mu^*, \sigma)}{\prod_i \text{dnorm}(y_i, \mu^{(s)}, \sigma)} \right) \left( \frac{\text{dnorm}(\mu^*, \mu_0, \sigma_0)}{\text{dnorm}(\mu^{(s)}, \mu_0, \sigma_0)} \right) \left( \frac{\text{dgamma}(\mu^*, \mu^{(s)}, 1)}{\text{dgamma}(\mu^{(s)}, \mu^*, 1)} \right). \quad (21)$$

# M-H example: Normal-Normal model with known $\sigma$ cont'd

$$\alpha^* = \left( \frac{\prod_i \text{dnorm}(y_i, \mu^*, \sigma)}{\prod_i \text{dnorm}(y_i, \mu^{(s)}, \sigma)} \right) \left( \frac{\text{dnorm}(\mu^*, \mu_0, \sigma_0)}{\text{dnorm}(\mu^{(s)}, \mu_0, \sigma_0)} \right) \left( \frac{\text{dgamma}(\mu^*, \mu^{(s)}, 1)}{\text{dgamma}(\mu^{(s)}, \mu^*, 1)} \right). \quad (22)$$

In many cases, computing the ratio  $r$  directly can be numerically unstable. Therefore, one can work with  $\log(r)$ .

$$\begin{aligned} \log(r) = & \sum_i \left( \log \text{dnorm}(y_i, \mu^*, \sigma) - \log \text{dnorm}(y_i, \mu^{(s)}, \sigma) \right) + \\ & \left( \log \text{dnorm}(\mu^*, \mu_0, \sigma_0) - \log \text{dnorm}(\mu^{(s)}, \mu_0, \sigma_0) \right) + \\ & \left( \log \text{dgamma}(\mu^*, \mu^{(s)}, 1) - \log \text{dgamma}(\mu^{(s)}, \mu^*, 1) \right). \end{aligned} \quad (23)$$

# M-H example: Normal-Normal model with known $\sigma$ cont'd

This is the Metropolis algorithm. How to update it for an M-H algorithm?

```
llik = sum(dnorm(y, mu, sigma, log = TRUE));

for (t in (thin + 1):iter){
  mup = runif(1, mu - C, mu + C);
  llikp = sum(dnorm(y, mup, sigma, log = TRUE));

  logr = llikp - llik + dnorm(mup, mu0, sigma0, log = TRUE) -
        dnorm(mu, mu0, sigma0, log = TRUE);

  logu = log(runif(1));
  if(logr > logu){
    mu = mup;
    llik = llikp;
    acc0 = acc0 + 1;
  }
```

# Outline

- 1 Overview
- 2 Metropolis Algorithm
- 3 Metropolis-Hastings Algorithm
- 4 Summary**



# Multi-parameter MCMC

With multiple parameters, a common strategy is to set up an MCMC sampler overall, and update each parameter using

- Draws from the full conditional when they are readily available (i.e. a Gibbs step).
- Draws from a Metropolis/M-H step otherwise.