

Lecture 9A

MTH6102: Bayesian Statistical Methods

Eftychia Solea

Queen Mary University of London

2023

Today's agenda

Today's lecture

- Understand Metropolis-Hastings
- Apply Metropolis-Hastings in Bayesian inference to generate samples from the posterior pdf.

Markov Chain Monte Carlo (MCMC)

- Recall, Monte Carlo integration approximates integrals of various functions $h(x)$

$$I = \int h(x)f(x) dx = E_f[h(X)], \quad X \sim f$$

by **directly sampling iid samples** from the pdf f or from the posterior pdf in Bayesian inference.

- Let X_1, \dots, X_n iid $\sim f$, the Monte Carlo estimator of I is given by

$$\hat{I} = \frac{1}{N} \sum_{i=1}^N h(X_i).$$

Markov Chain Monte Carlo (MCMC)

- **Question:** But what if we cannot sample directly from f ?
 - f is not analytically tractable.
 - Then, simple Monte Carlo integration cannot be used.
- In Bayesian inference, if we use a non-conjugate prior, then the posterior distribution may not be a well-known distribution.
 - our prior beliefs may not be captured using a conjugate prior
 - conjugate prior is unavailable for complicated problems

Motivating example

- Let $x = (x_1, \dots, x_n)$ IID from $N(\mu, \sigma^2)$, with μ known and σ^2 unknown.
- We showed that a gamma(α, β) prior for $\tau = 1/\sigma^2$ is conjugate.
- But what if a gamma(α, β) does not adequately represent our prior beliefs?

Motivating example

- Instead, we assume that our prior beliefs are represented by the **lognormal** (θ, v^2) distribution with pdf

$$p(\tau) = \frac{1}{\tau v \sqrt{2\pi}} \exp \left\{ -\frac{(\log \tau - \theta)^2}{2v^2} \right\}, \quad \tau > 0,$$

where θ and v^2 are known.

- What is the posterior density of τ under the lognormal prior and normal likelihood? What is the posterior mean of τ ?

Example

$x_1, \dots, x_n \sim \mathcal{N}(\mu, \sigma^2)$, where σ^2 is unknown and μ is known.

The likelihood is

$$p(x_1, \dots, x_n | \sigma^2) = \left(\frac{1}{2\pi}\right)^{n/2} T^{n/2} \exp\left\{-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2\right\}$$

The prior $p(\tau)$ is lognormal (θ, ν^2) with pdf

$$p(\tau) = \frac{1}{\tau \nu \sqrt{2\pi}} \exp\left\{-\frac{(\log \tau - \theta)^2}{2\nu^2}\right\}, \tau \in (0, \infty)$$

The posterior density, $p(\tau | x_1, \dots, x_n)$, is

$$p(\tau | x_1, \dots, x_n) = \frac{\text{prior} \times \text{likelihood}}{p(\text{data})}$$

$$= \frac{\text{Bayes numerator}}{p(\text{data})}$$

$$\frac{T^{n/2-1} \exp \left\{ -\frac{T}{2} \sum_{i=1}^n (\pi_i - \mu)^2 - \frac{(\log T - \theta)^2}{2\nu^2} \right\}}{\int_0^{\infty} T^{n/2-1} \exp \left\{ -\frac{T}{2} \sum_{i=1}^n (\pi_i - \mu)^2 - \frac{(\log T - \theta)^2}{2\nu^2} \right\} dT}$$

This is not a Gamma density

The normalising constant cannot be determined analytically
 We know the posterior up to a constant.

MCMC can help when f is not analytically tractable

- **Markov Chain Monte Carlo (MCMC)** is a set of methods that can generate a sample with pdf f without having to sample from f directly.
- Thus, MCMC can be used to generate samples from complicated probability distributions.
- At the price, however, of yielding **dependent** observations that are **approximately** from f .

Markov Chain Monte Carlo (MCMC)

- The general **idea** of MCMC methods is to construct a sequence of RV X_1, X_2, \dots , called **Markov chain**, which (hopefully) converges to the distribution of interest f .
- However, X_1, X_2, \dots , is NOT independent any more.
- But it can still be used to estimate means, $E_f[h(X)]$, because there is a WLLN for Markov chains.
- Under certain conditions,

$$\hat{I} = \frac{1}{N} \sum_{i=1}^N h(X_i) \xrightarrow{P} E_f[h(X)] = I, \quad \text{as } N \rightarrow \infty.$$

What is a Markov Chain?

- **Definition (Markov Chain).** A Markov chain is a sequence X_1, X_2, \dots of random variables such that the probability distribution of X_i (pmf or pdf) only depends on the previous value X_{i-1}

$$p(X_i | X_1, X_2, \dots, X_{i-2}, X_{i-1}) = p(X_i | X_{i-1}).$$

- The process depends on the past only through the present.

Example: Random walk

- As an example of a Markov chain is the **random walk** starting at $X_1 = 1$.
- Suppose $X_1 = 1$, and for $i > 1$

$$P(X_i = X_{i-1} + 1) = 1/2,$$

$$P(X_i = X_{i-1} - 1) = 1/2.$$

- So you flip a coin move +1 steps if heads, move -1 steps if tails.
- At step i of this Markov chain, X_{i-1} is either increased or decreased by 1 with probability $\frac{1}{2}$.

Metropolis-Hastings algorithm

- The **Metropolis-Hastings algorithm** is a type of MCMC that works as follows.
- Let $q(y|x)$ be a conditional density that we know how to sample from.
- $q(y|x)$ is called the **proposal distribution**.
- The Metropolis-Hastings algorithm creates a Markov Chain (dependent observations) X_1, X_2, \dots as follows.

Metropolis-Hastings algorithm

Goal: Generate X_1, X_2, \dots from f (target distribution)

Choose X_1 arbitrarily. Suppose we have generated X_1, \dots, X_i . To generate X_{i+1} do the following:

- 1 Generate a proposal or candidate random value $Y \sim q(y|X_i)$. (proposal distribution)
- 2 Evaluate $r \equiv r(X_i, Y)$ where

$$r(x, y) = \min \left\{ \frac{f(y) q(x|y)}{f(x) q(y|x)}, 1 \right\}.$$

- 3 Generate $U \sim U(0, 1)$. If $U < r$, set $X_{i+1} = Y$, otherwise set $X_{i+1} = X_i$.

Metropolis algorithm terminology

- q is the **proposal distribution**: we propose new rv Y using the conditional distribution $q(\cdot | X_i)$ that depends on X_i (not on the past).
- MH accepts Y with probability

$$r \equiv r(X_i, Y) = \min \left\{ \frac{f(Y) q(X_i|Y)}{f(X_i) q(Y|X_i)}, 1 \right\},$$

called the **acceptance probability**.

Metropolis algorithm terminology

- f is sometimes called the **target distribution**: this is what we are aiming for, i.e. we want to generate a sample with pdf f .
- In Bayesian inference, f would be the posterior distribution $p(\theta | y)$, and we want a sample of θ values from this posterior distribution.

Choosing events in computer code

Remarks:

- In general, to implement a random event that happens with probability r :
- Generate $u \sim \text{Uniform}(0, 1)$;
- Event happens if $u \leq r$.
- If U is a random variable, with $U \sim \text{Uniform}(0, 1)$, then U has cdf $F(r) = r$, so $P(U \leq r) = r$.

Metropolis-Hastings algorithm

Remarks:

- A common choice for $q(y|x)$ is $N(x, b^2)$ for some $b > 0$.
- This means that the proposal Y is drawn from normal centered at the current value.
- By symmetry, $q(y|x) = q(x|y)$

$$r(x, y) = \min \left\{ \frac{f(y)}{f(x)}, 1 \right\}.$$

$$q(y|x) \sim N(x, b^2) \quad b > 0$$

$$\text{Then } y \sim q(y|x_i) \sim N(x_i, b^2)$$

By symmetry, $q(y|x) = q(x|y)$.

Because, $q(y|x) \sim N(x, b^2)$

$$q(y|x) = \frac{1}{\sqrt{2\pi} b^2} \exp \left\{ -\frac{(y-x)^2}{2b^2} \right\}$$

$$= \frac{1}{\sqrt{2\pi} b^2} \exp \left\{ -\frac{(x-y)^2}{2b^2} \right\}$$

$$= q(x|y)$$

Then, my acceptance probability is

$$\min \left\{ \frac{f(y) q(x|y)}{f(x) q(y|x)}, 1 \right\} = \min \left\{ \frac{f(y)}{f(x)}, 1 \right\}$$

when $q(y|x)$ is symmetric we have the symmetric MH

Metropolis-Hastings algorithm

Remarks:

- In the algorithm, f only appears in acceptance probability

$$r(X_i, Y) = \min \left\{ 1, \frac{f(Y)}{f(X_i)} \right\}.$$

- The acceptance probability **does not depend on the normalisation constant**, i.e. if $f(x) = cg(x)$, where $c > 0$ doesn't depend on x , then

$$r(X_i, Y) = \min \left\{ 1, \frac{g(Y)}{g(X_i)} \right\}.$$

- So we only need to know f up to a normalisation constant. Useful for Bayesian inference!

. If $f(x) = cg(x)$, then

$$\min \left\{ \frac{f(y)}{f(x)}, 1 \right\} = \min \left\{ \frac{cg(y)}{cg(x)}, 1 \right\}$$

For Bayesian inference,

$$p(\theta|y) = \frac{p(\theta) \times \text{likelihood}}{P(\text{data})}$$

$$= C p(\theta) \times \text{likelihood}$$

Output of the Metropolis-Hastings algorithm

① The Metropolis-Hastings algorithm generates a dependent sequence of observations X_1, X_2, \dots

② Since our procedure for generating X_{i+1} depends only on X_i , the conditional distribution of X_{i+1} given X_1, \dots, X_i depends only on X_i .

③ Hence, the sequence X_1, X_2, \dots is a Markov chain.

$$y \sim q(y|X_i)$$

Output of the Metropolis-Hastings algorithm

- The chain X_1, X_2, \dots has the property that:
if $X_{i-1} \sim f$, then $X_i \sim f$.
- f is the equilibrium distribution or stationary of the chain.
- However, we don't start with $X_1 \sim f$ (because if we could, we wouldn't need this algorithm).
- But for large enough i , if some technical conditions are met, then each $X_i \sim f$ approximately.

$(X_i)_i$ is an approximate sample from f

Output of the Metropolis-Hastings algorithm

- 1 In practice, we only generate X_1, X_2, \dots, X_N for some large N .
- 2 Under some conditions, the empirical distribution of X_1, X_2, \dots, X_N approximates f well if N is large.
- 3 Hence, we can approximate the integral $I = \int h(x)f(x) dx$ using the approximated X_1, X_2, \dots, X_N , that is

$$\hat{I} = \frac{1}{N} \sum_{i=1}^N h(X_i), \quad X_1, X_2, \dots, X_N \sim f \text{ (approximately),}$$

and X_1, X_2, \dots, X_N generated by MH.

Example: Metropolis-Hastings algorithm

- The Cauchy distribution has density

$$f(x) = \frac{1}{\pi(1+x^2)}$$

target distribution

- Our goal is to simulate a Markov chain whose stationary distribution is f .
- Take $q(y|x)$ to be $N(x, b^2)$ for some $b > 0$.
- Then,

$$r(x, y) = \min \left\{ \frac{1+x^2}{1+y^2}, 1 \right\}.$$

- Let $r = r(X_i, Y)$. Generate $U \sim U(0, 1)$. If $U < r$, set $X_{i+1} = Y$, otherwise set $X_{i+1} = X_i$.

Goal: Generate a sample $\{X_1, \dots, X_n\}$ from f , where

$$f(x) = \frac{1}{n(1+x^2)} \quad (x \in \mathbb{R})$$

Set X_1 randomly. Suppose we have generated $\{X_1, \dots, X_i\}$. We add a new observation X_{i+1} by doing the following

① Generate $Y \sim N(X_i, b^2)$, where b is fixed.

② Compute

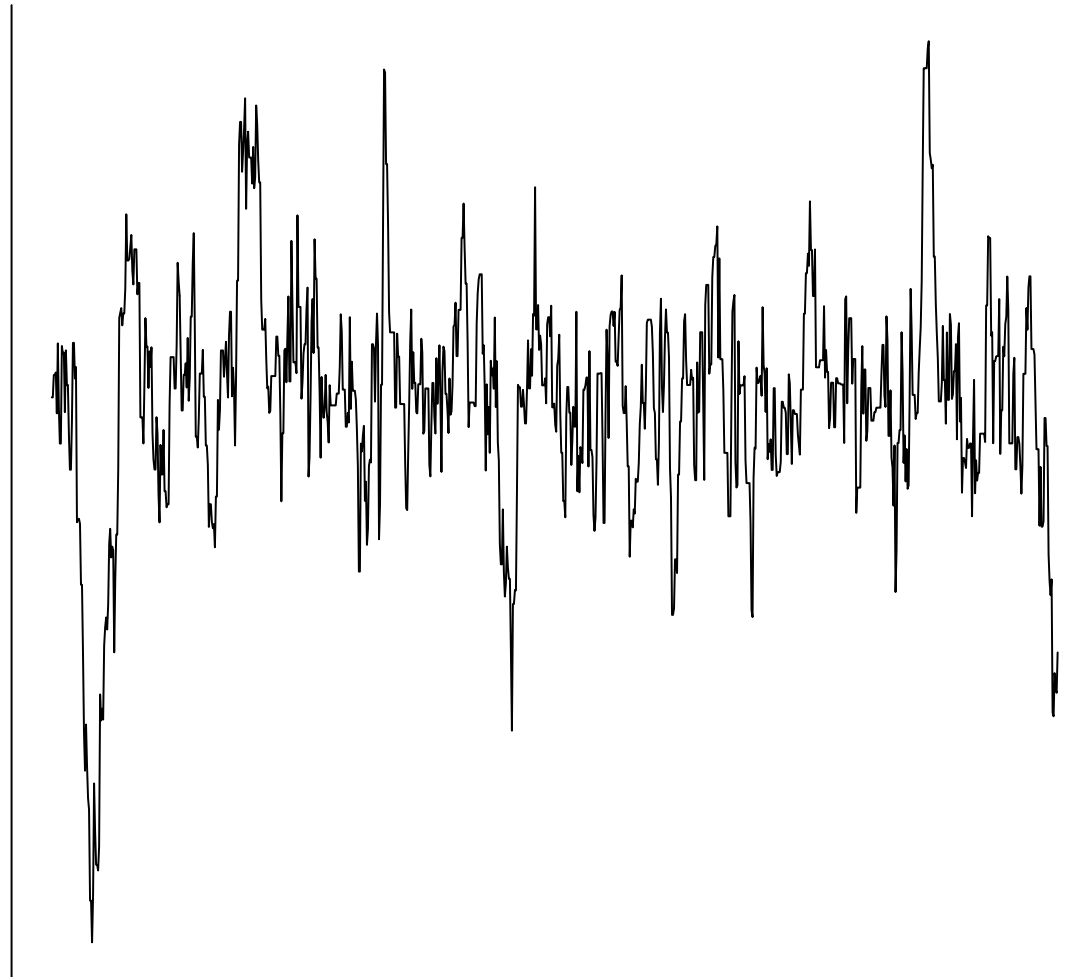
$$\begin{aligned} r \equiv r(X_i, Y) &= \min \left\{ \frac{f(Y)}{f(X_i)}, 1 \right\} \\ &= \min \left\{ \frac{\frac{1}{n(1+Y^2)}}{\frac{1}{n(1+X_i^2)}}, 1 \right\} \end{aligned}$$

$$= \min \left\{ \frac{1+X_i^2}{1+Y^2}, 1 \right\} \quad (\text{acceptance probability})$$

③ Generate $U \sim U[0,1]$
if $U < r$, then $X_{i+1} = Y$
otherwise $X_{i+1} = X_i$

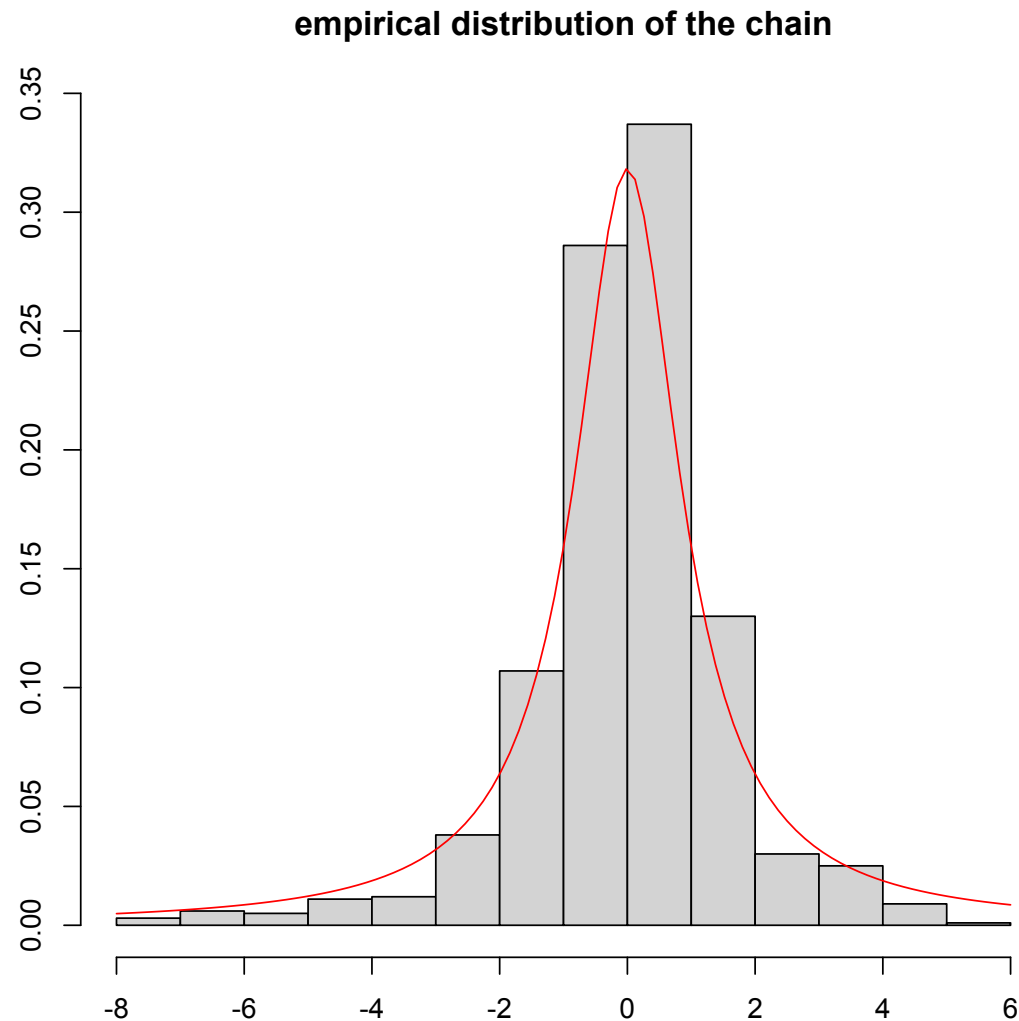
Example: Metropolis-Hastings algorithm

- Figure below shows the chains of length $N = 1000$ using $b = 1$



Example: Metropolis-Hastings algorithm

- Figure: Histogram of chains and the plot of the Cauchy density (red)
- The distribution of chain converges to the desired Cauchy distribution.



Relevance to Bayesian inference

- Let f be the posterior pdf, $p(\theta | y)$: this is the distribution we want to sample from.
- Let $q(\psi | \theta_i)$ be a pdf for the proposal ψ which is symmetric in ψ and θ , e.g., normal $N(\theta_i, b^2)$.
 $\psi \sim N(\theta_i, b^2)$
- The algorithm constructs a Markov chain $\theta_1, \theta_2, \dots$, where the θ_i are continuous rvs (in our applications).

Relevance to Bayesian inference

- q is called the proposal distribution: it is used to generate the next possible point in the Markov chain.
- q is often taken as a normal distribution centred on the current point

$$\psi_i \sim N(\theta_i, b^2), \text{ for some } b > 0.$$

- The normal pdf is symmetric in θ and ψ , as required by the algorithm

$$q(\psi | \theta) = \frac{1}{\sqrt{2\pi}b} e^{-\frac{(\psi - \theta)^2}{2b^2}} = \frac{1}{\sqrt{2\pi}b} e^{-\frac{(\theta - \psi)^2}{2b^2}} = q(\theta | \psi).$$

Relevance to Bayesian inference

Goal: Generate a posterior sample $\{\theta_1, \theta_2, \dots\}$ from $p(\theta|y)$

The algorithm constructs a Markov chain $\theta_1, \theta_2, \dots$ as follows:

- Start with arbitrary θ_1 .

- For each $i > 1$, generate ψ from distribution $q(\psi | \theta_i)$.

- Let

$$r = \min \left\{ 1, \frac{p(\psi | y)}{p(\theta_i | y)} \right\}$$

- Set

$$\theta_{i+1} = \begin{cases} \psi & \text{with probability } r \\ \theta_i & \text{with probability } 1 - r \end{cases}$$

target distribution

$= q(\theta_i | \psi)$
e.g. $N(\theta_i | \sigma^2)$

Review

Relevance to Bayesian inference

- In Bayesian inference, the posterior density is

$$p(\theta | y) \propto \underbrace{p(\theta)} \underbrace{p(y | \theta)}$$

- It's difficult to find the normalizing constant

$$\int \underbrace{p(\theta) p(y | \theta)} d\theta$$

- We don't need to find this, we just put $g(\theta) = \underbrace{p(\theta) p(y | \theta)}$, use g in the algorithm (where we have f), and we will get an approximate sample from $p(\theta | y)$.
- The Markov chain $\theta_1, \theta_2, \dots$ is this sample.

$$p(\theta|y) = \frac{p(\theta) \times p(y|\theta)}{\int p(\theta) \times p(y|\theta) d\theta} = c p(\theta) \times p(y|\theta),$$

where c is the normalising constant.
 We don't need to compute c to implement MH to compute the acceptance probability, because

$$\begin{aligned} r &= \min \left\{ \frac{p(\psi|y)}{p(\theta_i|y)}, 1 \right\} \\ &= \min \left\{ \frac{p(\psi) p(y|\psi)}{p(\theta_i) p(y|\theta_i)}, 1 \right\} \\ &= \min \left\{ \frac{p(\psi) p(y|\psi)}{p(\theta_i) p(y|\theta_i)}, 1 \right\} \end{aligned}$$

We only need to know the Bayes numerator
 $p(\theta) \times p(y|\theta)$

Metropolis algorithm for Bayesian inference

Define $g(\theta) = p(\theta) p(y | \theta)$, the non-normalized posterior density.

Generate a Markov chain $\theta_1, \theta_2, \dots$ as follows:

- Choose some $b > 0$.
- Start with θ_1 , where $g(\theta_1) > 0$.
- For each $i > 1$:

- Generate $\psi \sim N(\theta_i, b^2)$.

- Let

$$r = \min \left\{ 1, \frac{g(\psi)}{g(\theta_i)} \right\}.$$

- Set

$$\theta_{i+1} = \begin{cases} \psi & \text{with probability } r \\ \theta_i & \text{with probability } 1 - r \end{cases}$$

Working on the log scale

- We usually do the computations using the log of the posterior density.
- The likelihood is typically a product of many terms.

$$p(y | \theta) = \prod_{i=1}^n p(y_i | \theta)$$

- Due to finite accuracy of computers, if we multiply these together for a large dataset, the result is inaccurate.
- So calculate

$$\log(p(y | \theta)) = \sum_{i=1}^n \log(p(y_i | \theta))$$

Using the log scale

$$= \log p(\theta | y)$$

- Define $\mathcal{L}(\theta) = \log(p(\theta) p(y | \theta)) = \log(p(\theta)) + \log(p(y | \theta))$, the log of the posterior density (up to a constant).
- To work on the log scale, the part of the algorithm with the acceptance probability changes.

- Define

$$\delta = \min(0, \mathcal{L}(\psi) - \mathcal{L}(\theta_{i-1}))$$

- Generate $u \sim \text{Uniform}(0, 1)$
- Set

$$\theta_i = \begin{cases} \psi & \text{if } \log(u) < \delta \\ \theta_{i-1} & \text{otherwise} \end{cases}$$

$$r = \min \left\{ \frac{p(\psi|y)}{p(\theta_i|y)}, 1 \right\}$$

At the log-scale

$$\begin{aligned} \log \left(\frac{p(\psi|y)}{p(\theta_i|y)} \right) &= \log p(\psi|y) - \log p(\theta_i|y) \\ &= \log(p(\psi) p(y|\psi)) - \log[p(\theta_i) p(y|\theta_i)] \\ &= L(\psi) - L(\theta_i) \end{aligned}$$

At the log-scale

$$r = \min \left\{ L(\psi) - L(\theta_i), 0 \right\}$$

$$U \sim U[0,1]$$

$$\text{if } \log U < r$$

$$\theta_{i+1} = \psi$$

$$\text{else } \theta_{i+1} = \theta_i$$

Normal example with known variance

- Y_1, \dots, Y_n iid from $N(\theta, \sigma^2)$ where σ^2 is known.
- $\theta \sim N(\mu, \tau^2)$ with τ^2 known,
- Apply the Metropolis-Hastings algorithm to simulate from the posterior $p(\theta|y_1, \dots, y_n)$ after observing $Y = y$

Metropolis-Hastings algorithm for Bayesian inference

- Metropolis-Hastings algorithm generates a dependent sequence $\theta^{(1)}, \dots$, of θ values.
- Under mild conditions, the empirical distribution of $\theta^{(i)}$, $i = 1, 2, \dots$ will approximate well the posterior.
- We can view $\theta^{(i)}$, $i = 1, 2, \dots$ as a sample from the posterior $p(\theta|y)$.
- Hence, we can approximate posterior means, quantiles and other posterior quantities of interest using $\{\theta^{(1)}, \dots, \theta^{(N)}\}$ for large N .
- However, our approximation to these quantities will depend on how well our simulated sequence actually approximates $p(\theta | y)$.