# MARKOV CHAIN MONTE CARLO

BEE 6940 LECTURE 8                              MARCH 13, 2023

# TABLE OF CONTENTS

# REVIEW OF MARKOV CHAINS

# Markov Chains

A Markov chain is a stochastic process based on memoryless probabilistic transitions between states.

**Key properties**:

- Detailed Balance (reversibility) $\Rightarrow$ Existence of stationary distribution

- Ergodicity (irreducible and aperiodic) $\Rightarrow$ Stationary distribution is unique and can be obtained as a limiting distribution.

# IDEA OF SAMPLING ALGORITHM

If we construct an ergodic Markov chain with the appropriate transition probabilities for detailed balance to hold with respect to the target distribution, we can use the chain realizations as samples from the target distribution.

- No need for closed-form representation of the posterior

- Can use these samples as normal (with some caveats) for means, quantiles, simulations, etc.

This category of sampling algorithms is called **Markov chain Monte Carlo (MCMC)**.

# CONVERGENCE TO THE TARGET DISTRIBUTION

Given a Markov chain $\{X_t\}_{t=1,\ldots,T}$ returned from this procedure, sampling from distribution $\pi$:

- $\mathbb{P}(X_t = y) \to \pi(y)$ as $t \to \infty$

- This means the chain can be considered a *dependent* sample approximately distributed from $\pi$.

- The first values (the *transient portion*) of the chain are highly dependent on the initial value.

# The Metropolis-Hastings Algorithm

# METROPOLIS-HASTINGS ALGORITHM

The **Metropolis-Hastings** algorithm is the foundational MCMC algorithm (and was named one of the top ten algorithms of the 20th century).

It builds a Markov chain based on transitions by accepting/rejecting proposals for new samples from a *conditional proposal distribution $q(y|x)$* and accepting or rejecting those proposals.

# METROPOLIS-HASTINGS ALGORITHM

Given $X_t = x_t$:

1. Generate $Y_t \sim q(y|x_t)$;

2. Set $X_{t+1} = Y_t$ with probability $\rho(x_t, Y_t)$, where

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

else set $X_{t+1} = x_t$.

# Metropolis-Hastings Algorithm

That almost seems too simple!

But the devil is in the details: performance and efficiency are highly dependent on the choice of $q$.

# METROPOLIS-HASTINGS ALGORITHM

That almost seems too simple!

But the devil is in the details: performance and efficiency are highly dependent on the choice of $q$.

**Key**: There is a tradeoff between exploration and acceptance.

- Wide proposal: Can make bigger jumps, may be more likely to reject proposals.

- Narrow proposal: More likely to accept proposals, may not "mix" efficiently.

# M-H and Detailed Balance

Why do M-H transitions satisfy detailed balance?

Write detailed balance equation as:

$$\pi(x)k(x, y) = \pi(y)k(y, x),$$

where $\pi(\cdot)$ is the stationary distribution and $k(x, y)$ is the *transition kernel* (not necessarily normalized) of moving from $x \to y$.

# M-H AND DETAILED BALANCE

The transition kernel for M-H transitions can be written as:

$$k(x, y) = q(y|x)\rho(x, y),$$

where $q$ is the proposal distribution and $\rho$ is the M-H acceptance probability.

# M-H AND DETAILED BALANCE

If the proposal is rejected, *e.g.* $x = y$, then detailed balance is satisfied trivially.

Else, the left-hand side of the detailed balance equation becomes:

$$\pi(x)k(x,y) = \pi(x)q(y|x)\rho(x,y).$$

.

# M-H and Detailed Balance

Without loss of generality, suppose $\pi(y)q(x|y) < \pi(x)q(y|x)$, so $\rho(x,y) < 1$. Then

$$
\begin{aligned}
\pi(y)k(x,y) &= \pi(x)q(y|x) \times \frac{\pi(y)}{\pi(x)}\frac{q(x|y)}{q(y|x)} \\
&= \pi(y)q(x|y) \\
&= \pi(y)k(y,x)
\end{aligned}
$$

since $\rho(y,x) = 1$.

# CONSIDERATIONS FOR IMPLEMENTATION

# PROPOSAL DISTRIBUTIONS

The original Metropolis et al (1953) algorithm focused on symmetric distributions $(q(y|x) = q(x|y))$, which are a convenient choice as then the acceptance probability reduces to

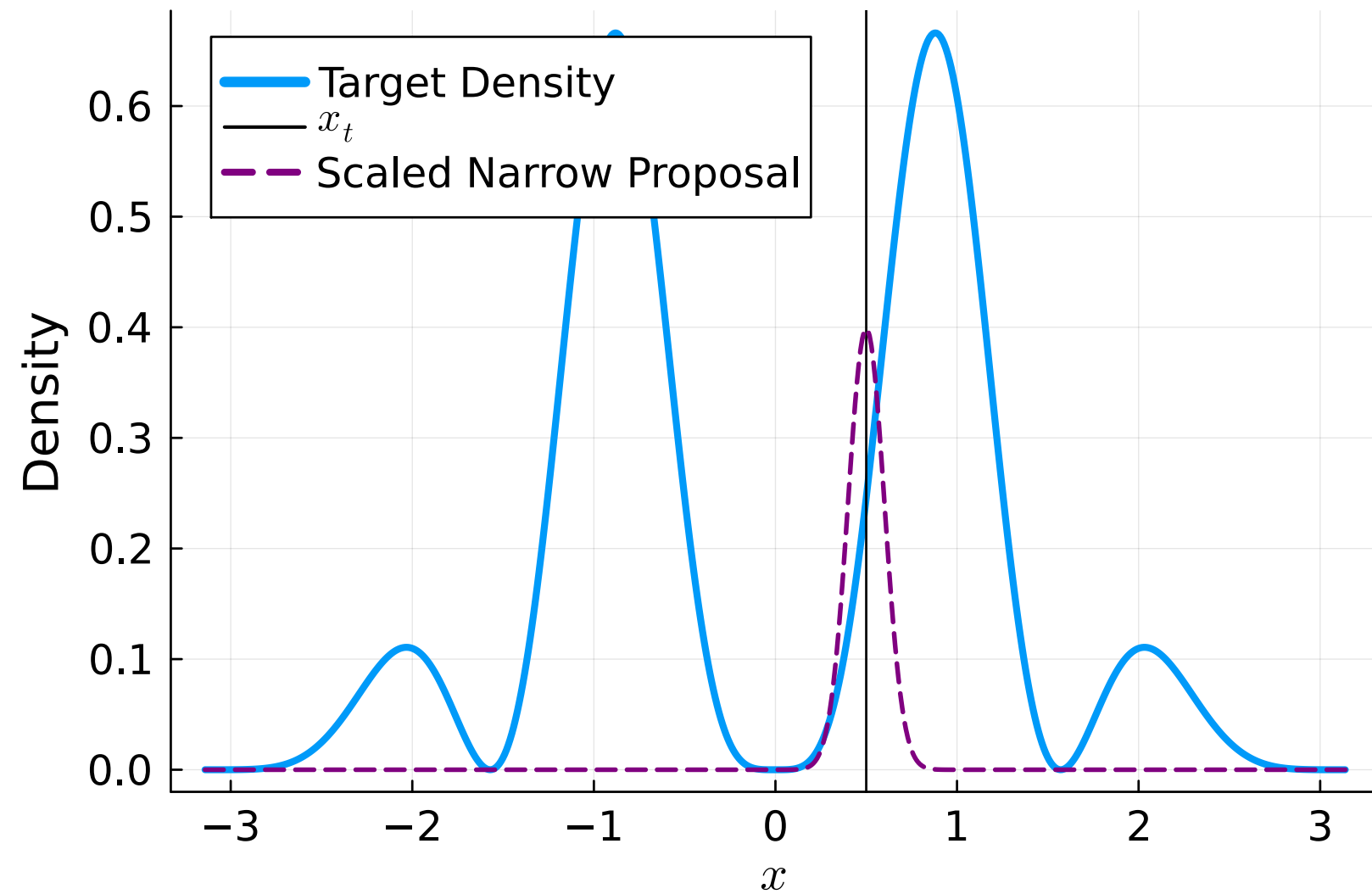$$\rho = \min\left\{\frac{\pi(y)}{\pi(x)}, 1\right\}.$$

# PROPOSAL DISTRIBUTIONS

For example, a common choice is a normal density $y \sim \mathrm{Normal}(x_t, \sigma^2)$ centered around the current point.

This turns the "exploration" part of the M-H algorithm into a random walk.

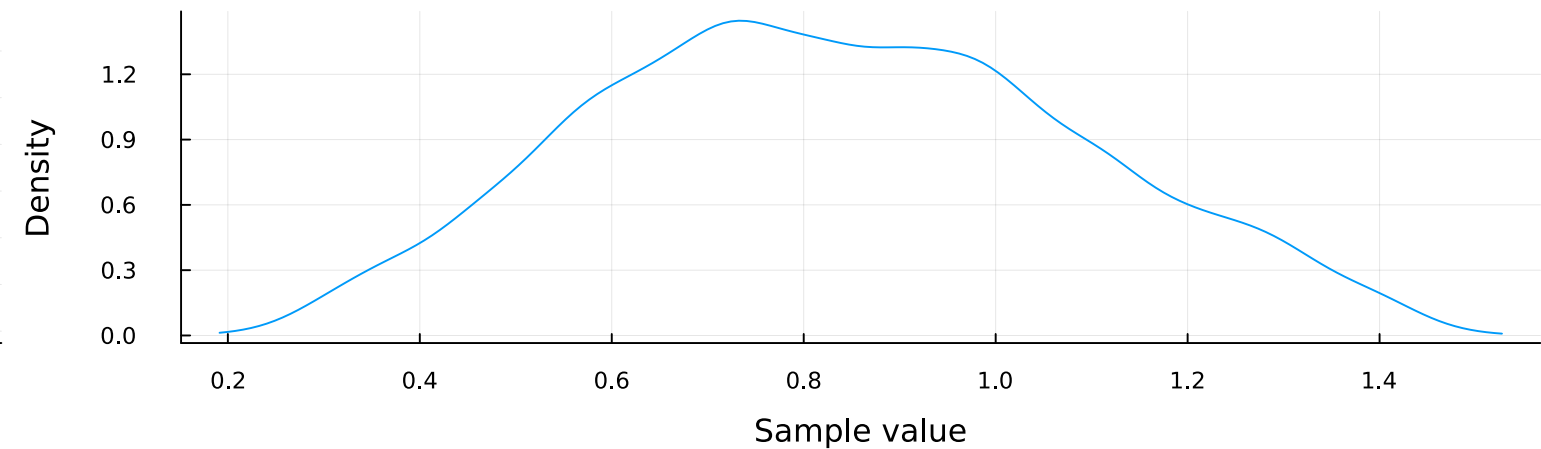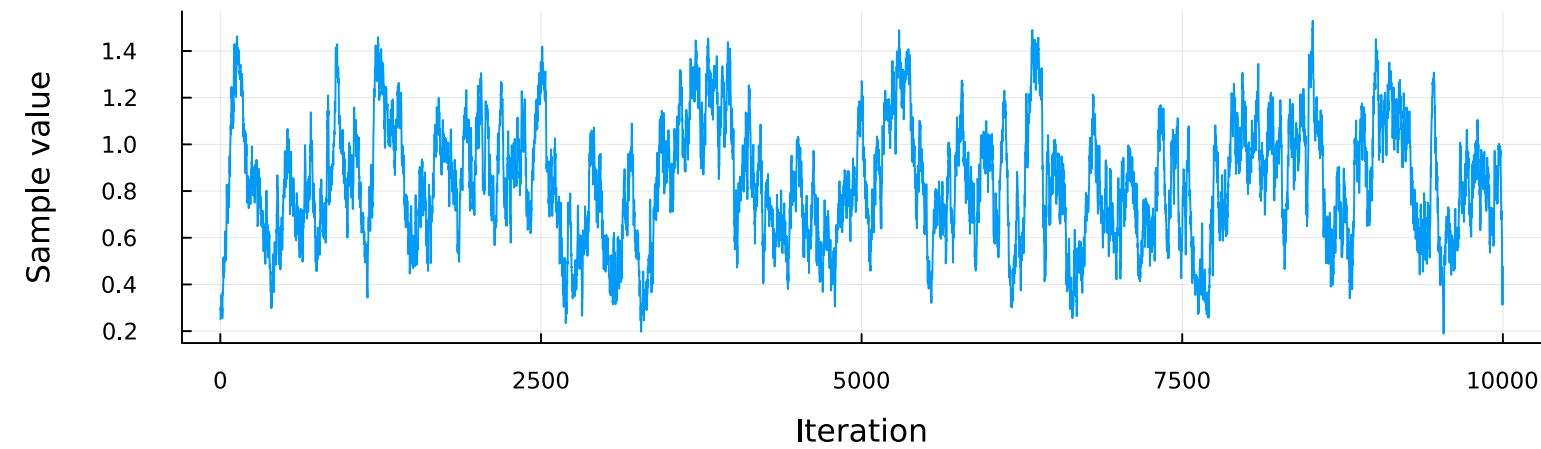# PROPOSAL DISTRIBUTIONS: EXAMPLE

# Sampling Efficiency

Two common measures of sampling efficiency:

- **Acceptance Rate**: Rate at which proposals are accepted

  - "Optimally" 30-45% (depending on number of parameters)

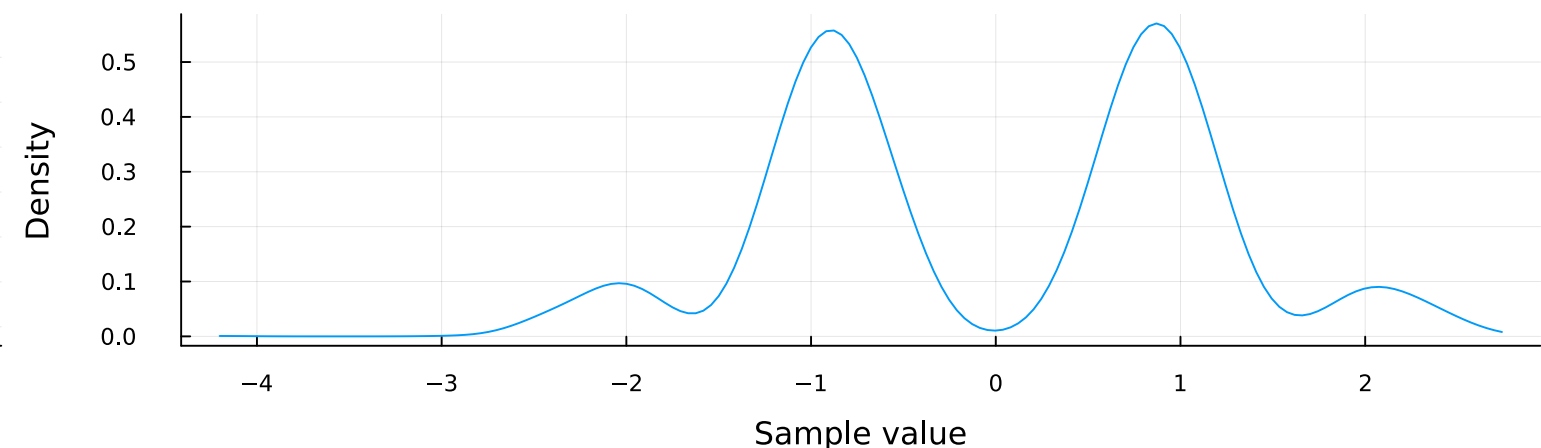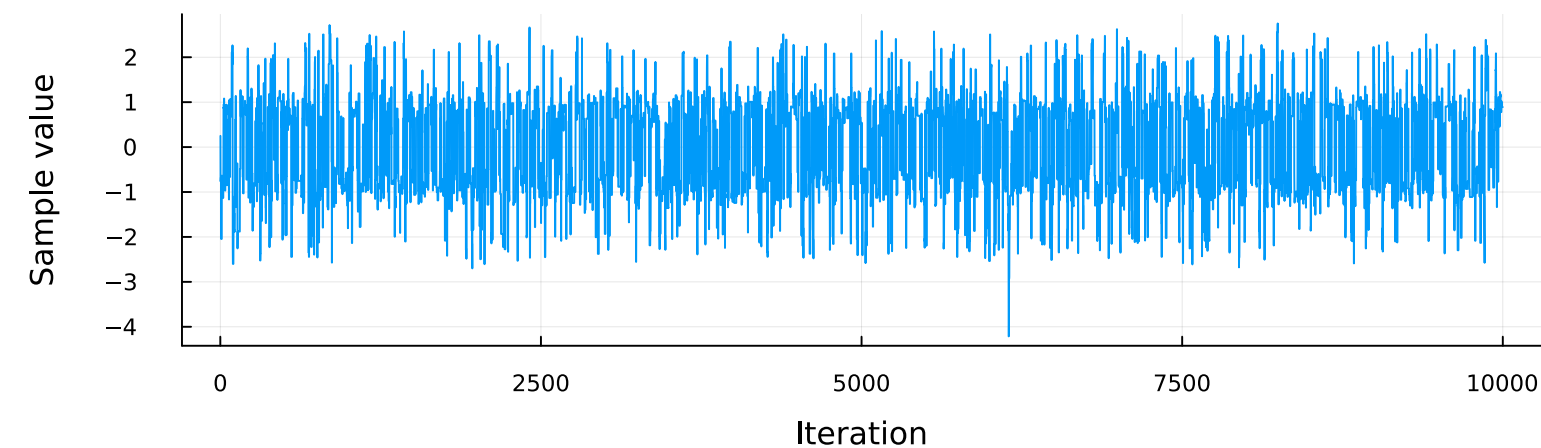- **Effective Sample Size (ESS)**: Accounts for autocorrelation $\rho_t$ across samples

$$N_{\text{eff}} = \frac{N}{1 + 2 \sum_{t=1}^{\infty} \rho_t}$$
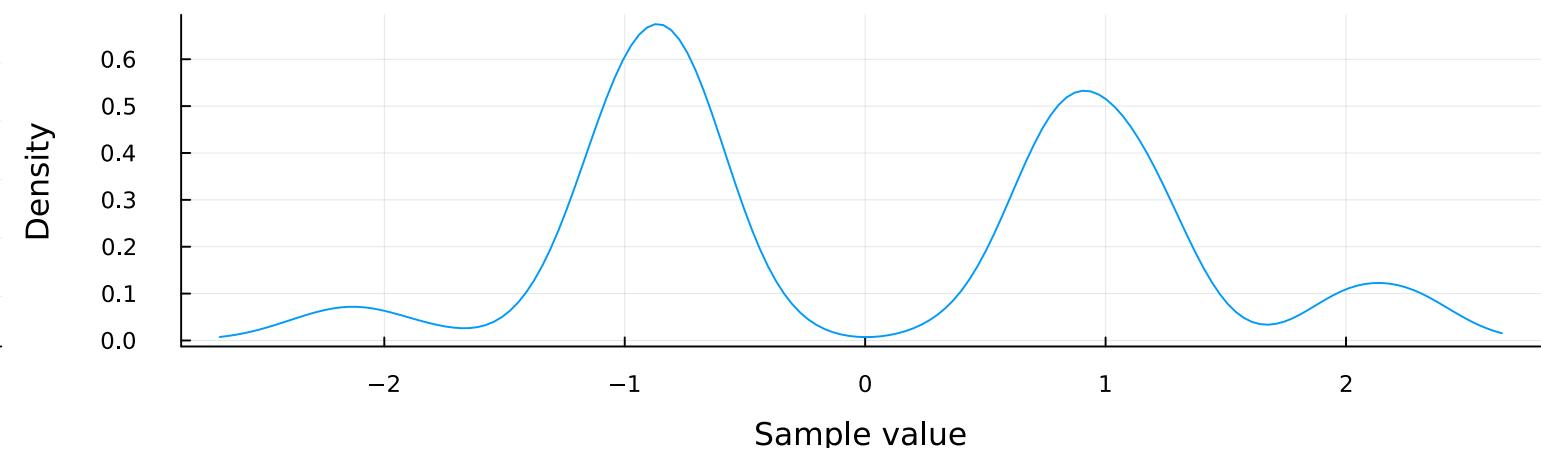
# Sampling Efficiency Example

# Autocorrelation Of Chains

# ESS By Proposal Variance

# ASSESSING CONVERGENCE

# CONVERGENCE TO STATIONARY DISTRIBUTION

Since the samples are a Markov chain, there is a *transient* portion prior to convergence to the stationary distribution.

What to do about these samples?

# CONVERGENCE TO STATIONARY DISTRIBUTION

Since the samples are a Markov chain, there is a *transient* portion prior to convergence to the stationary distribution.

What to do about these samples?

- Discard as *burn-in*;

- Just run the chain longer.

# HOW TO IDENTIFY CONVERGENCE?

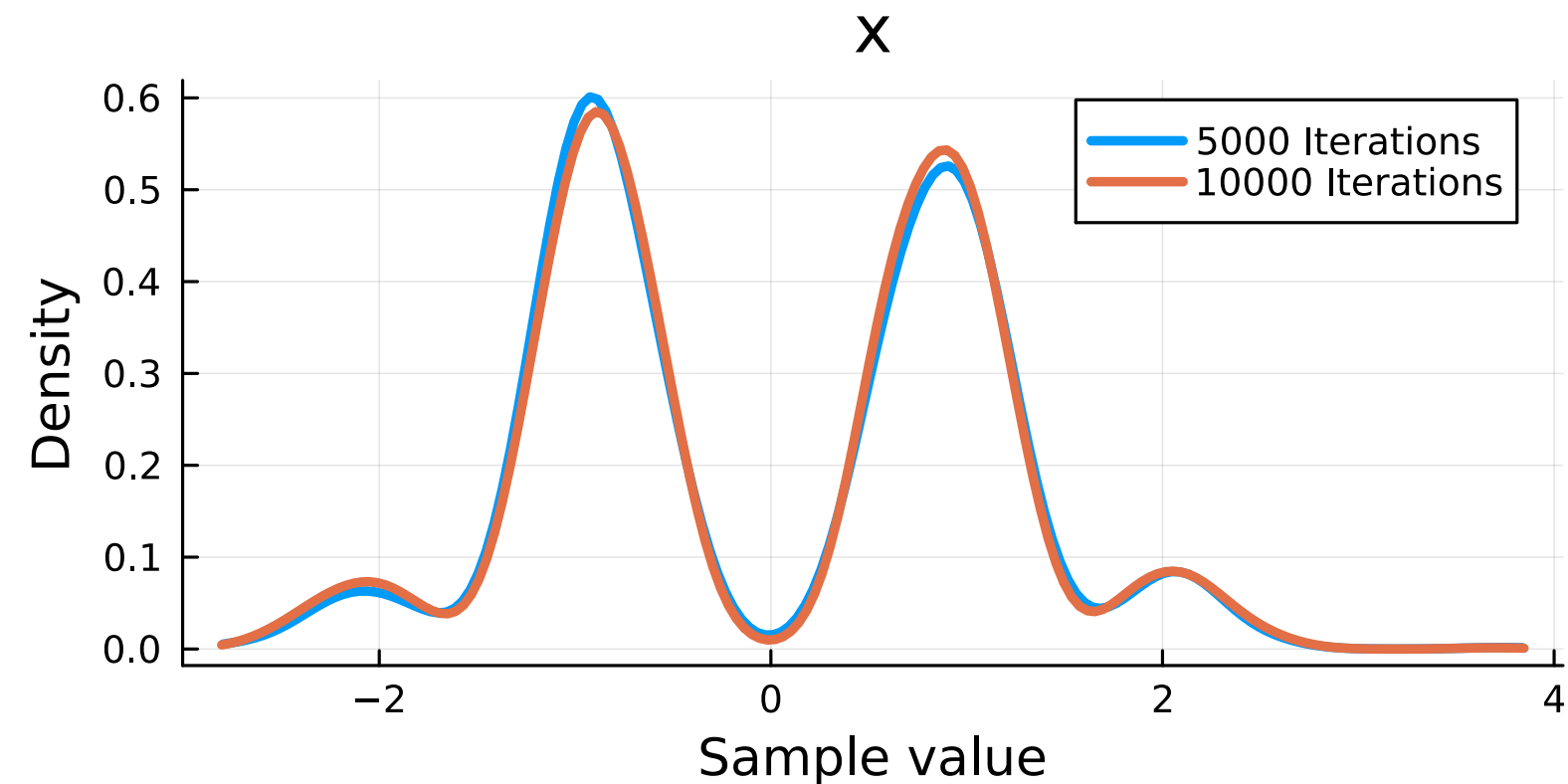This is probably the most challenging part of MCMC, other than tuning the proposal distribution.

**Short answer**: There is no guarantee! Judgement based on an accumulation of evidence from various heuristics.

The good news — getting the precise "right" end of the transient chain doesn't matter. If a few transient iterations of the transient portion, the effect will be washed out with a large enough post-convergence chain.

# HEURISTICS FOR CONVERGENCE

- Compare distribution (histogram/kernel density plot) after half of the chain to full chain.

# HEURISTICS FOR CONVERGENCE

- Gelman-Rubin criterion (Gelman & Rubin (1992)):

  - Run multiple chains from "overdispersed" starting points

  - Compare intra-chain and inter-chain variances

  - Summarized as $\hat{R}$ statistic: closer to 1 implies better convergence.

- Can also check distributions across multiple chains vs. the half-chain check. This is the default in `Turing.jl` with multiple chains (will see in lab).

# ASIDE: MULTIPLE CHAINS

Unless a specific scheme is used, multiple chains are not a solution for issues of convergence, as each individual chain needs to converge and have burn-in discarded/watered-down.

This means multiple chains are more useful for diagnostics, but once they've all been run long enough, can mix samples freely.

# Heuristics for Convergence

- If you're more interested in the mean estimate, can also look at the its stability by iteration or the *Monte Carlo standard error*.

- Look at traceplots; do you see sudden "jumps"?

- **When in doubt, run the chain longer.**

# APPROACHES FOR INCREASING EFFICIENCY

- Adaptive M-H (*e.g.* Vihola (2012)): adjusts proposal density to hit target acceptance rate

  - Need to be cautious about detailed balance

  - Typical strategy is to adapt for a portion of the initial chain (part of the burn-in), then run longer with that proposal.

# APPROACHES FOR INCREASING EFFICIENCY

- Hamiltonian Monte Carlo

    - Idea: Use proposals which steer towards "typical set" without collapsing towards the mode (based on Hamiltonian vector field);

    - Requires gradient information: can be obtained through autodifferentiation; challenging for external, black-box models;

    - Can be very efficient due to potential for anti-correlated samples, but efficiency is also be sensitive to parameterization.

# KEY TAKEAWAYS: MCMC

- Construct ergodic and reversible Markov chains with posterior as stationary distribution.

- Metropolis-Hastings: conceptually simple algorithm, but implementation plays a major role.

- Must rely on "accumulation of evidence" from heuristics for determination about convergence to stationary distribution.

- Transient portion of chain: Meh. Some people worry about this too much. Discard or run the chain longer.

- Parallelizing solves few problems, but running multiple chains can be useful for diagnostics.

# Upcoming Schedule

# Upcoming Schedule

**Wednesday**: Lab on using `Turing.jl` for probabilistic programming and MCMC

**Next Monday**: Storm surge and modeling extreme values