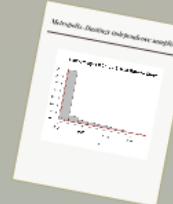
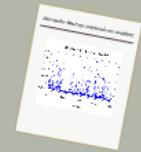
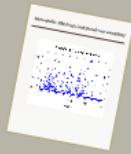
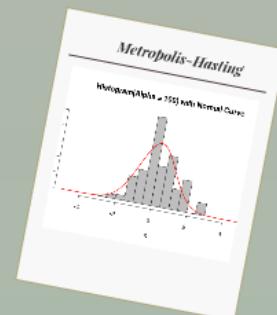
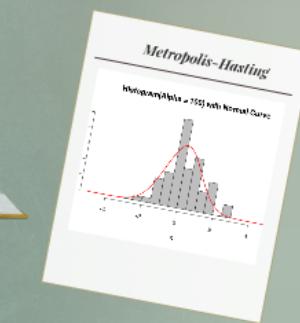
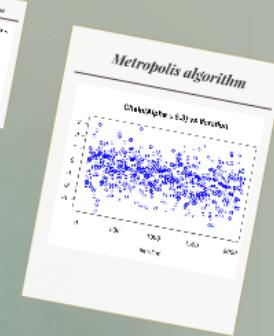
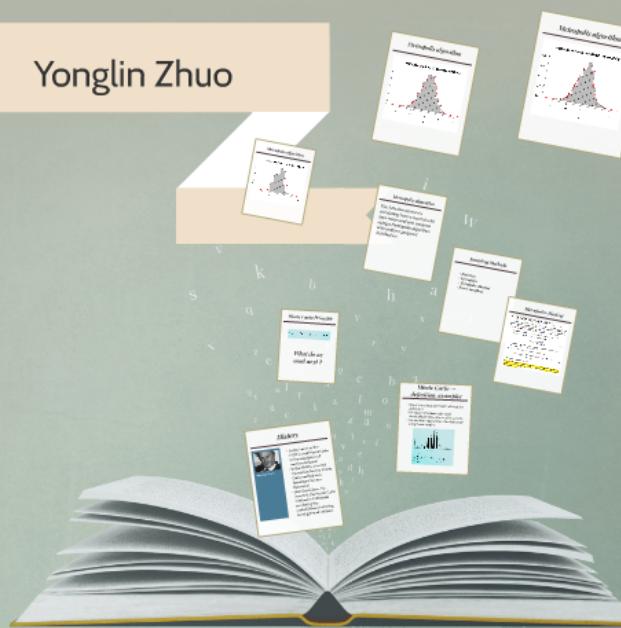


MCMC Yonglin Zhuo



MCMC

Yonglin Zhuo



History



Stanislaw Ulam

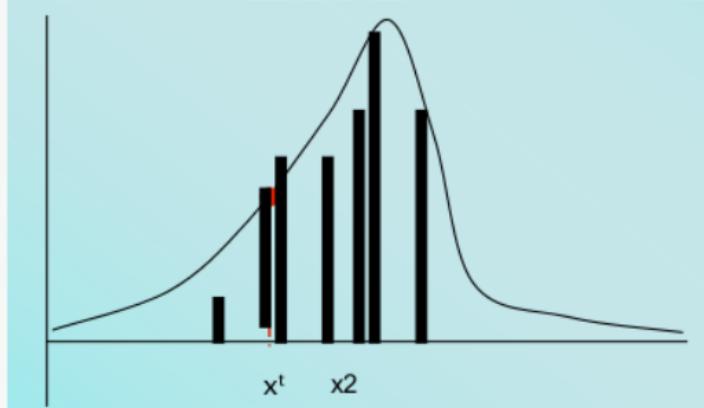
- Enrico Fermi in the 1930's used Monte Carlo in the calculation of neutron diffusion
- In the 1940's, a formal foundation for the Monte Carlo method was developed by von Neumann
- Stanislaw Ulam, He invented the Monte Carlo method in 1946 while pondering the probabilities of winning a card game of solitaire.

i

O

Monte Carlo -- definition, examples

- Given a very large set X and a distribution $p(x)$ over it
- We draw i.i.d. (independent and identically distributed) a set of N samples
- We can then approximate the distribution using these samples



$$\frac{1}{N} \sum_{i=1}^N \mathbf{1}(x^{(i)} = x) \xrightarrow{N \rightarrow \infty} p(x)$$

Monte Carlo Principle

$$E_N(f) = \frac{1}{N} \sum_{i=1}^N f(x^{(i)}) \xrightarrow{N \rightarrow \infty} E(f) = \sum_x f(x) p(x)$$

*What do we
need next ?*

Sampling Methods

- *Rejection*
- *Metropolis*
- *Metropolis-Hastings*
- *Exact Sampling*

a

j

X

Metrop

7.1 METROPOLIS-HASTINGS

A very general method for constructing a Markov chain Monte Carlo algorithm [324, 460]. The method begins by drawing a candidate point \mathbf{x}^* drawn at random from some starting distribution $f(\mathbf{x}^{(0)}) > 0$. Given $\mathbf{X}^{(t)} = \mathbf{x}^{(t)}$, the algorithm proceeds as follows:

1. Sample a candidate value \mathbf{X}^* from the proposal distribution $R(\mathbf{u}, \mathbf{x}^{(t)})$.
2. Compute the Metropolis-Hastings acceptance probability

$R(\mathbf{u}, \mathbf{x}^*)$

Note that $R(\mathbf{x}^{(t)}, \mathbf{X}^*)$ is always non-negative. This can only occur if $f(\mathbf{x}^{(t)}) > 0$ and $R(\mathbf{u}, \mathbf{x}^*) < 1$.

Metropolis-Hastings

7.1 METROPOLIS–HASTINGS ALGORITHM

A very general method for constructing a Markov chain is the Metropolis–Hastings algorithm [324, 460]. The method begins at $t = 0$ with the selection of $\mathbf{X}^{(0)} = \mathbf{x}^{(0)}$ drawn at random from some starting distribution g , with the requirement that $f(\mathbf{x}^{(0)}) > 0$. Given $\mathbf{X}^{(t)} = \mathbf{x}^{(t)}$, the algorithm generates $\mathbf{X}^{(t+1)}$ as follows:

1. Sample a candidate value \mathbf{X}^* from a *proposal distribution* $g(\cdot | \mathbf{x}^{(t)})$.
2. Compute the *Metropolis–Hastings ratio* $R(\mathbf{x}^{(t)}, \mathbf{X}^*)$, where

Metropolis-Hastings

7.1 METROPOLIS-HASTINGS ALGORITHM

A very general method for constructing a Markov chain is the Metropolis–Hastings algorithm [324, 460]. The method begins at $t = 0$ with the selection of $\mathbf{X}^{(0)} = \mathbf{x}^{(0)}$ drawn at random from some starting distribution g , with the requirement that $f(\mathbf{x}^{(0)}) > 0$. Given $\mathbf{X}^{(t)} = \mathbf{x}^{(t)}$, the algorithm generates $\mathbf{X}^{(t+1)}$ as follows:

1. Sample a candidate value \mathbf{X}^* from a *proposal distribution* $g(\cdot | \mathbf{x}^{(t)})$.
2. Compute the *Metropolis–Hastings ratio* $R(\mathbf{x}^{(t)}, \mathbf{X}^*)$, where

$$R(\mathbf{u}, \mathbf{v}) = \frac{f(\mathbf{v}) g(\mathbf{u} | \mathbf{v})}{f(\mathbf{u}) g(\mathbf{v} | \mathbf{u})}. \quad (7.1)$$

Note that $R(\mathbf{x}^{(t)}, \mathbf{X}^*)$ is always defined, because the proposal $\mathbf{X}^* = \mathbf{x}^*$ can only occur if $f(\mathbf{x}^{(t)}) > 0$ and $g(\mathbf{x}^* | \mathbf{x}^{(t)}) > 0$.

3. Sample a value for $\mathbf{X}^{(t+1)}$ according to the following:

$$\mathbf{X}^{(t+1)} = \begin{cases} \mathbf{X}^* & \text{with probability } \min\{R(\mathbf{x}^{(t)}, \mathbf{X}^*), 1\}, \\ \mathbf{x}^{(t)} & \text{otherwise.} \end{cases} \quad (7.2)$$

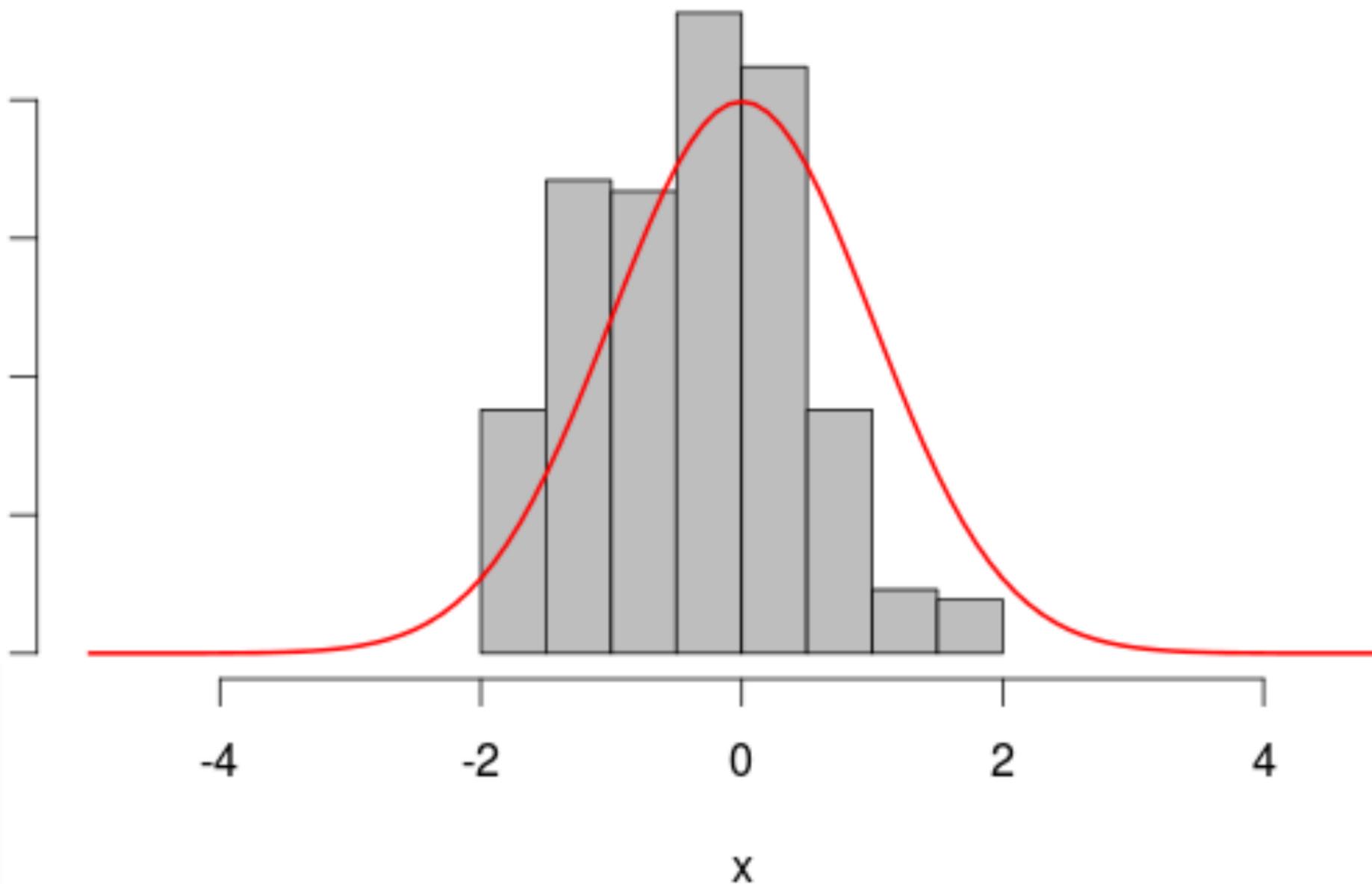
4. Increment t and return to step 1.

We will call the t th iteration the process that generates $\mathbf{X}^{(t)} = \mathbf{x}^{(t)}$. When the proposal distribution is symmetric so that $g(\mathbf{x}^{(t)} | \mathbf{x}^*) = g(\mathbf{x}^* | \mathbf{x}^{(t)})$, the method is known as the Metropolis algorithm [460].

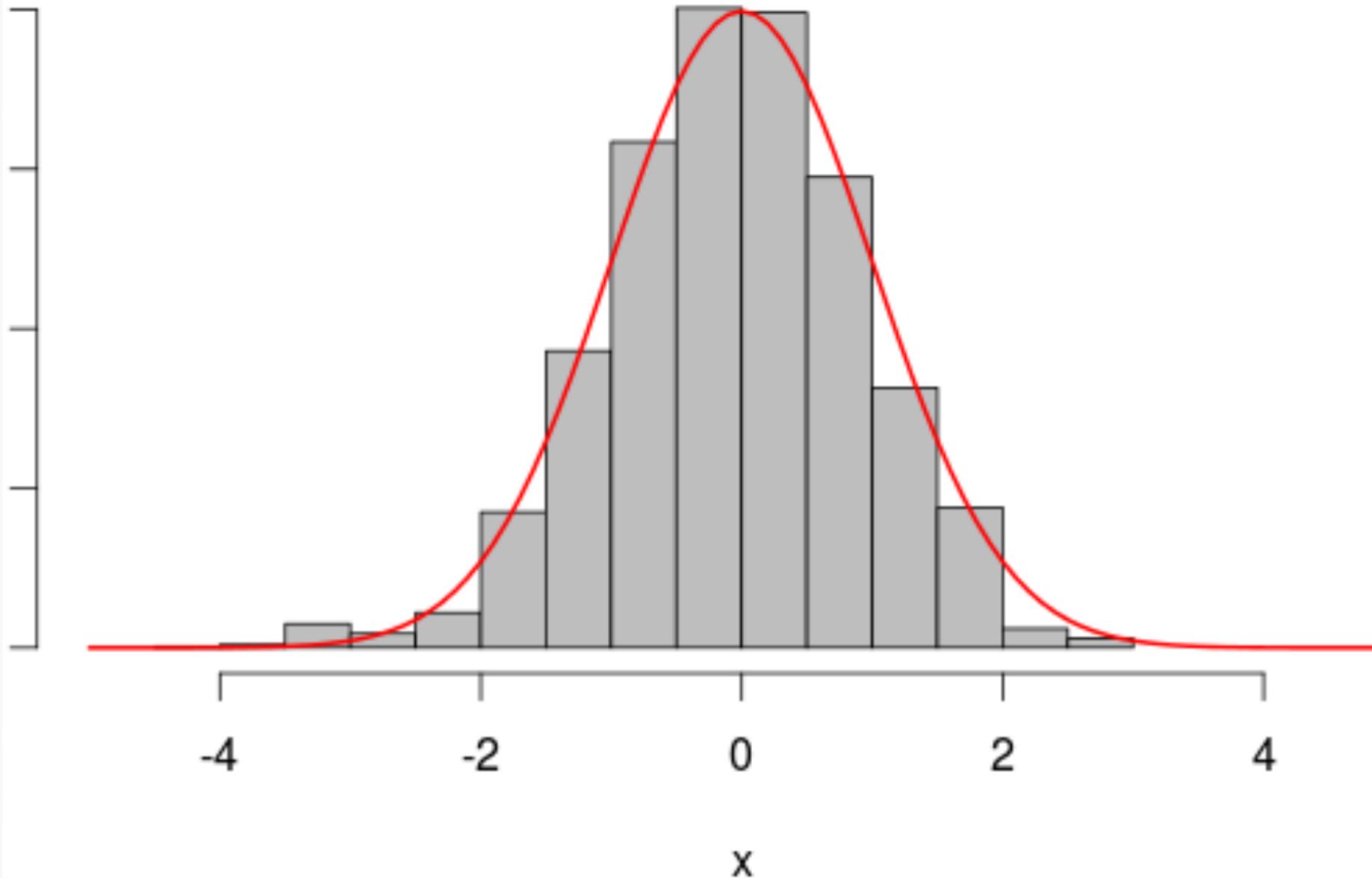
Metropolis algorithm

The function `normm` is simulating from a normal with zero mean and unit variance using a Metropolis algorithm with uniform proposal distribution.

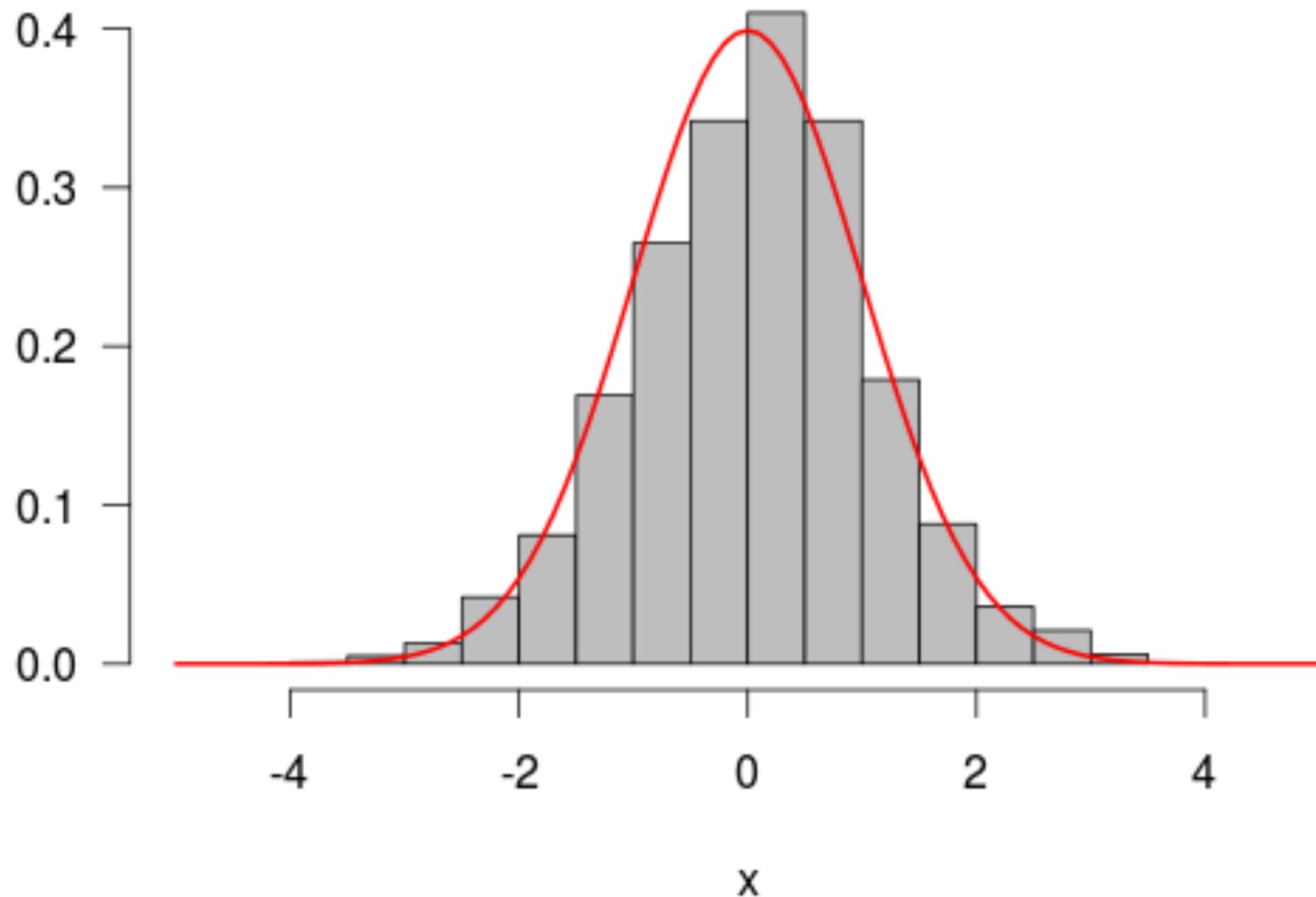
Histogram(Alpha = 0.1) with Normal Curve



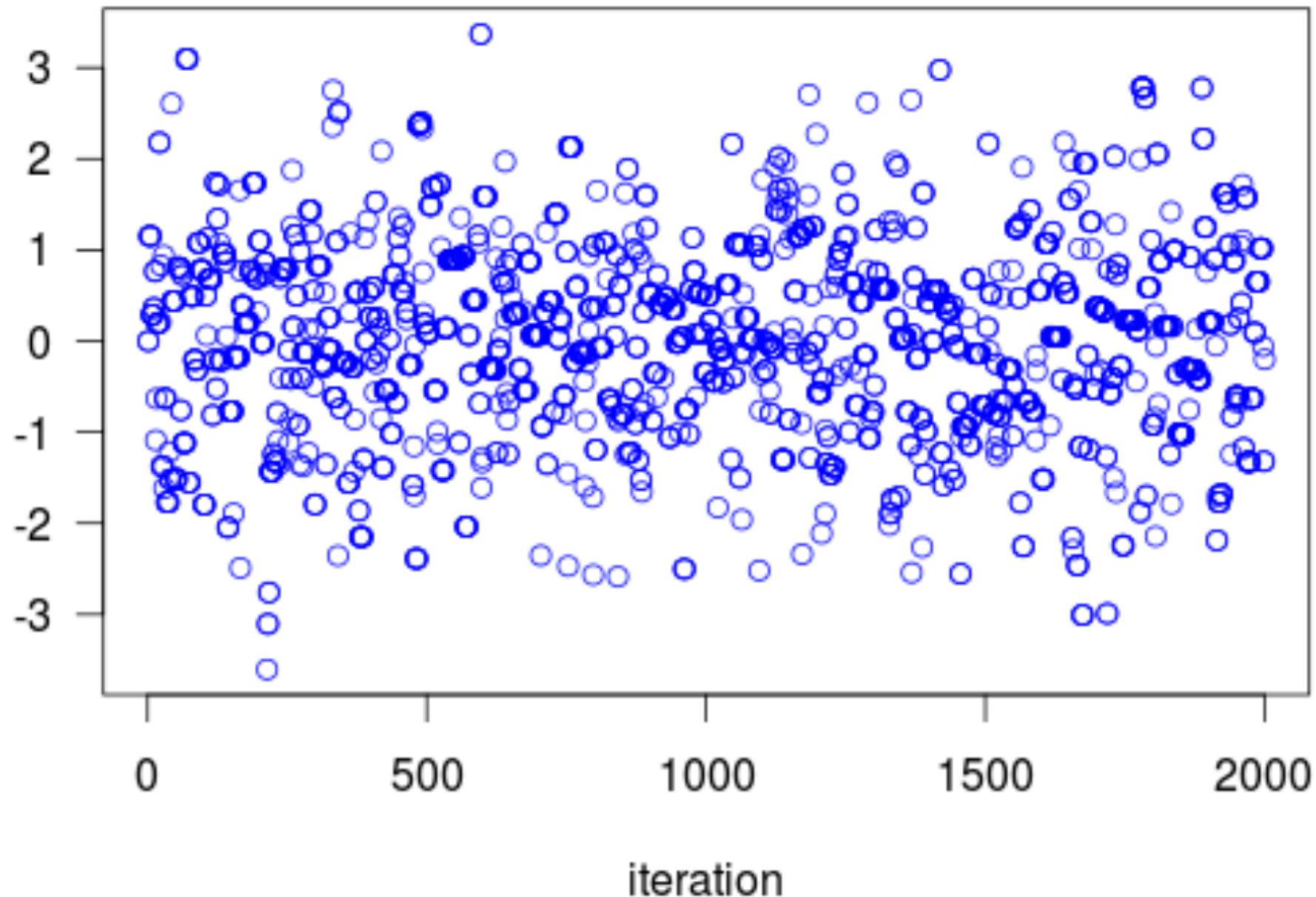
Histogram(Alpha = 1) with Normal Curve



Histogram(Alpha = 5.3) with Normal Curve



Chain(Alpha = 5.3) vs Iteration



Suppose that the proposal distribution for the Metropolis–Hastings algorithm is chosen such that $g(\mathbf{x}^* | \mathbf{x}^{(t)}) = g(\mathbf{x}^*)$ for some fixed density g . This yields an independence chain, where each candidate value is drawn independently of the past. In this case, the Metropolis–Hastings ratio is

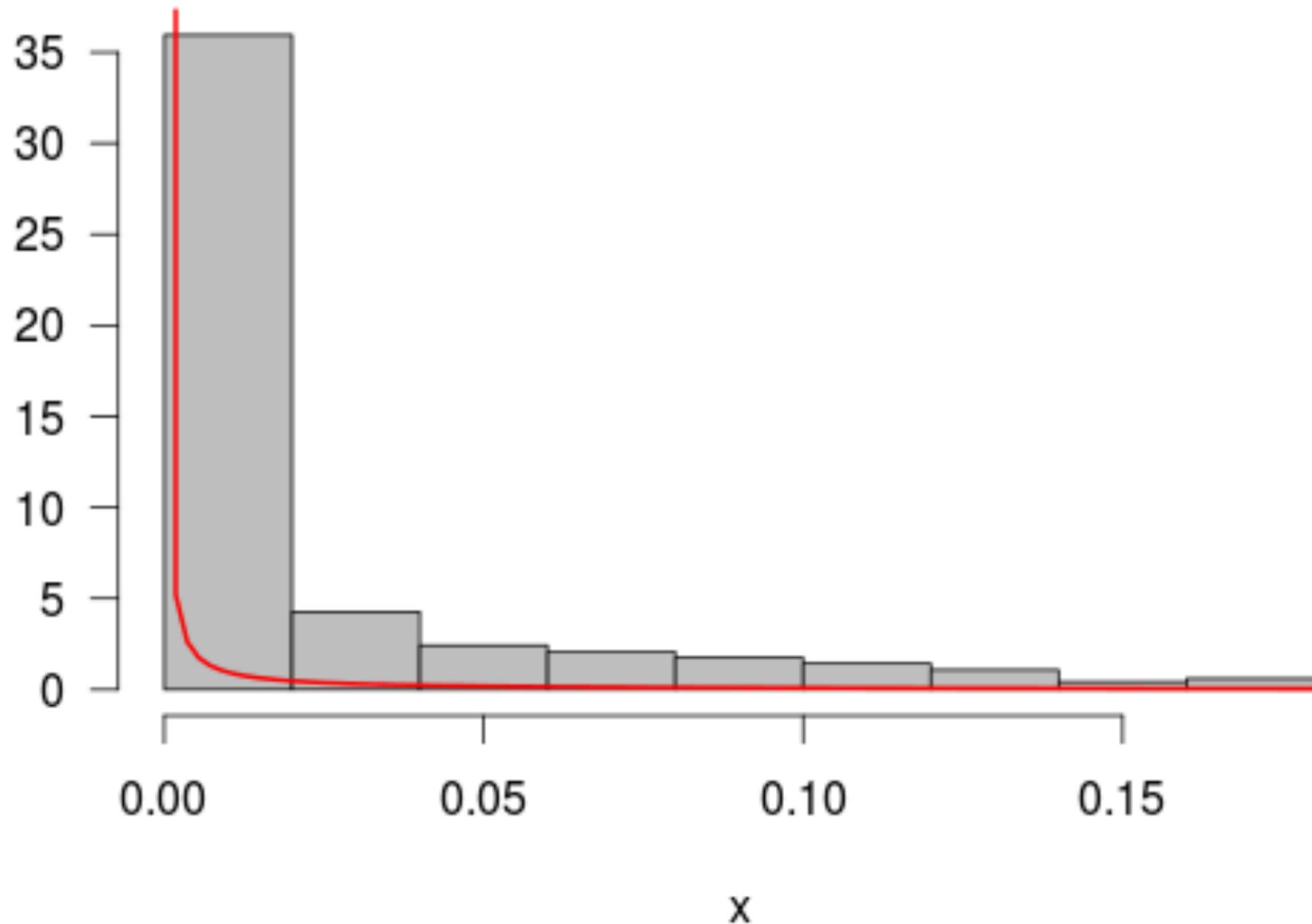
$$R(\mathbf{x}^{(t)}, \mathbf{X}^*) = \frac{f(\mathbf{X}^*) g(\mathbf{x}^{(t)})}{f(\mathbf{x}^{(t)}) g(\mathbf{X}^*)}. \quad (7.4)$$

The resulting Markov chain is irreducible and aperiodic if $g(\mathbf{x}) > 0$ whenever $f(\mathbf{x}) > 0$.

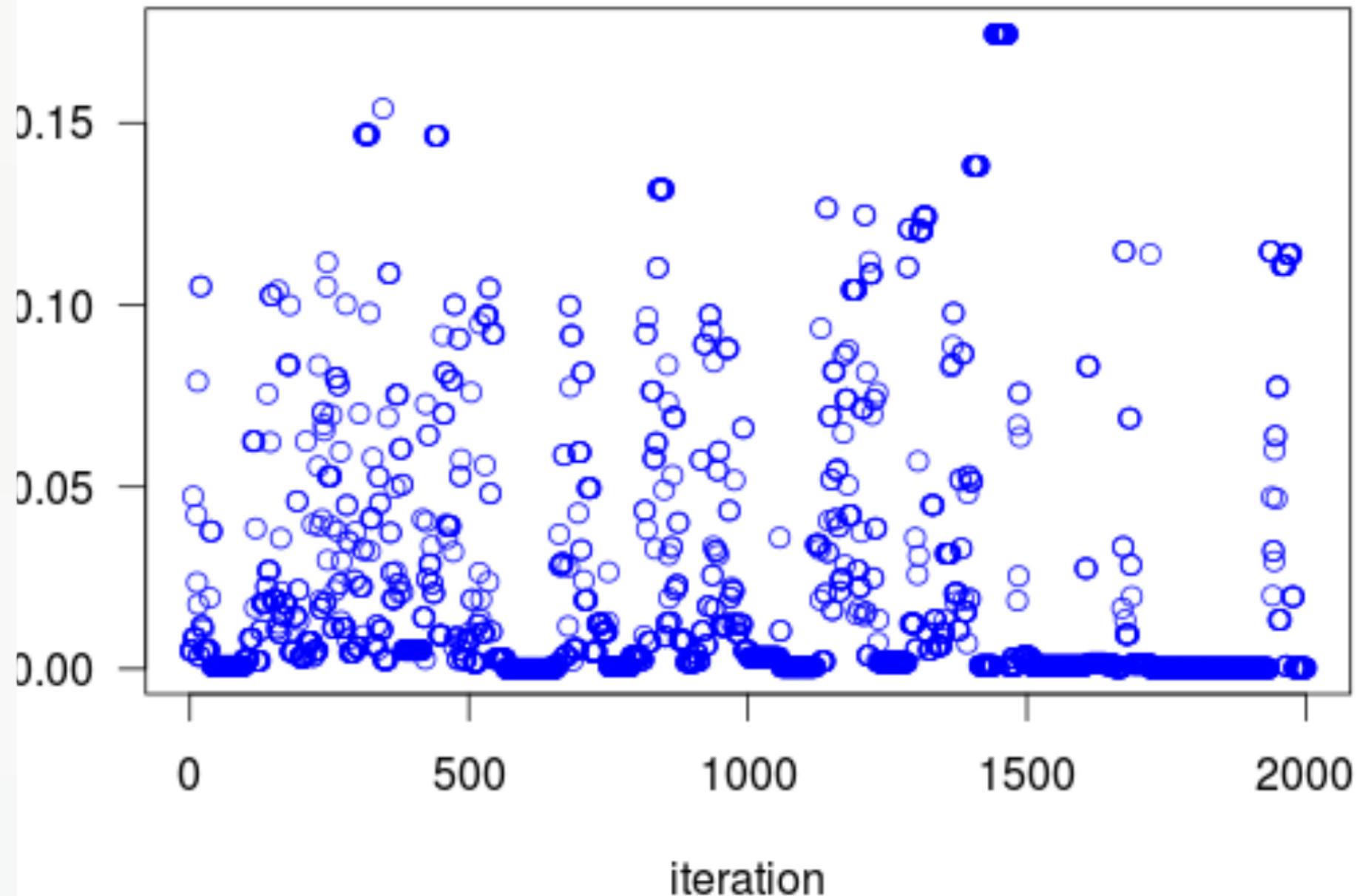
Notice that the Metropolis–Hastings ratio in (7.4) can be reexpressed as the ratio of importance ratios (see Section 6.4.1) where f is the target and g is the envelope: If $w^* = f(\mathbf{X}^*) / g(\mathbf{X}^*)$ and $w^{(t)} = f(\mathbf{x}^{(t)}) / g(\mathbf{x}^{(t)})$, then $R(\mathbf{x}^{(t)}, \mathbf{X}^*) = w^* / w^{(t)}$. This reexpression indicates that when $w^{(t)}$ is much larger than typical w^* values, then the chain will tend to get stuck for long periods at the current value. Therefore, the criteria discussed in Section 6.3.1 for choosing importance sampling envelopes are also relevant here for choosing proposal distributions: The proposal distribution g should resemble the target distribution f , but should cover f in the tails.

Example 7.1 (Bayesian Inference) MCMC methods like the Metropolis–Hastings algorithm are particularly popular tools for Bayesian inference, where some data \mathbf{y} are observed with likelihood function $L(\boldsymbol{\theta} | \mathbf{y})$ for parameters $\boldsymbol{\theta}$ which have prior distribution $p(\boldsymbol{\theta})$. Bayesian inference is based on the posterior distribution $p(\boldsymbol{\theta} | \mathbf{y}) =$

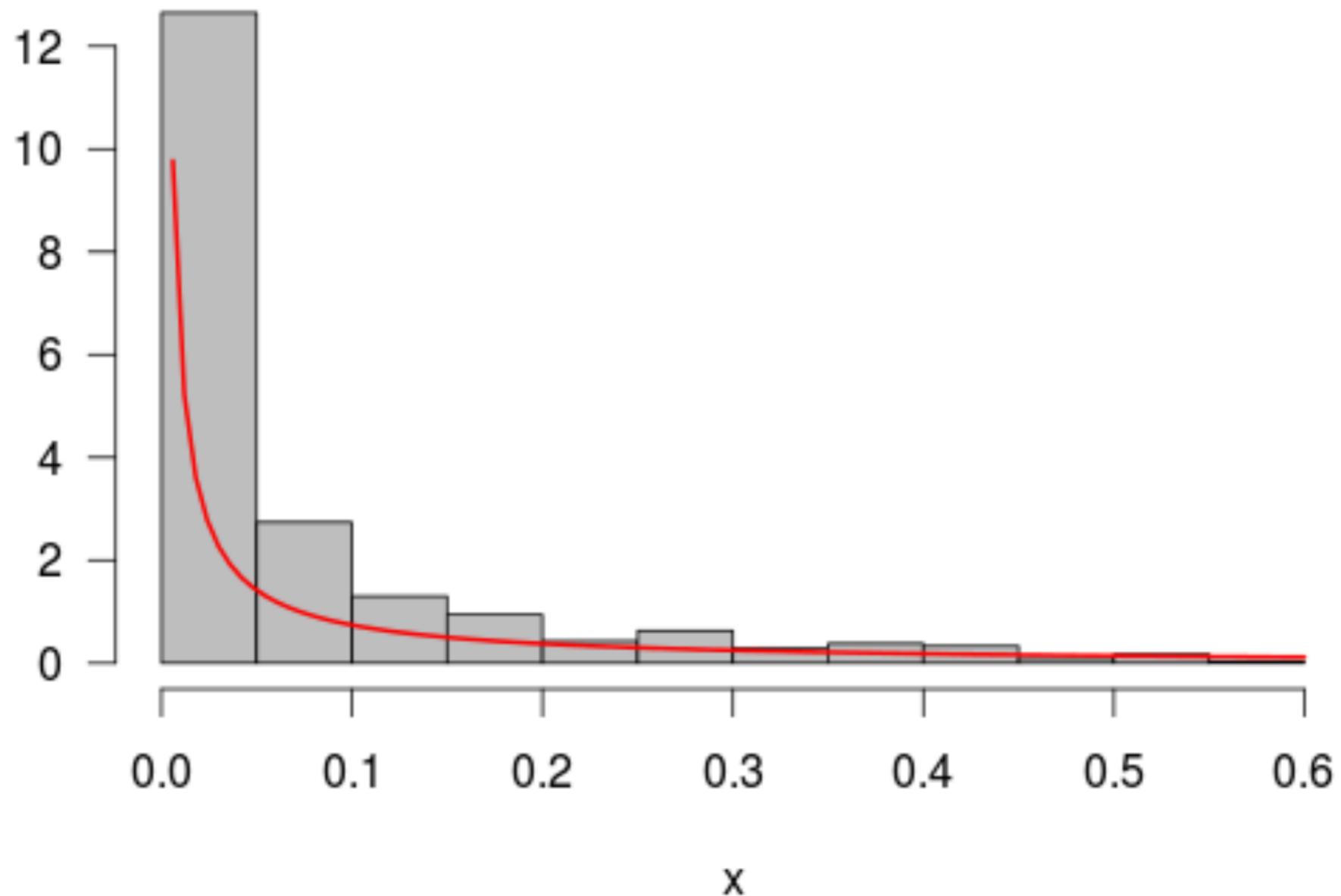
Histogram($a = 0.01$, $b = 2$) with Gamma Curve



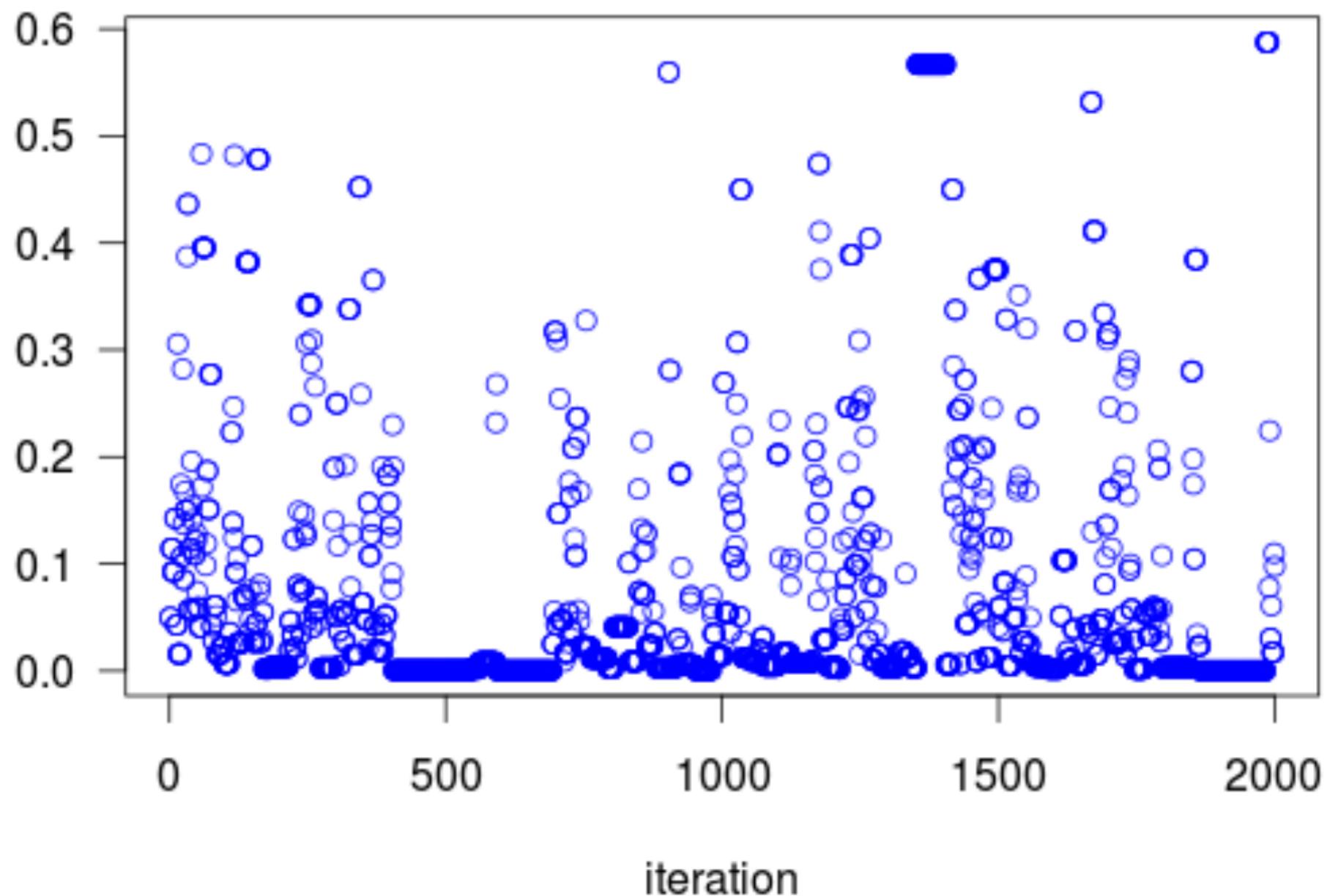
Chain($a = 0.01$, $b = 2$) vs Iteration



Histogram($a = 0.1$, $b = 2$) with Gamma Curve



Chain($a = 0.1$, $b = 2$) vs Iteration



MCMC

Yonglin Zhuo

