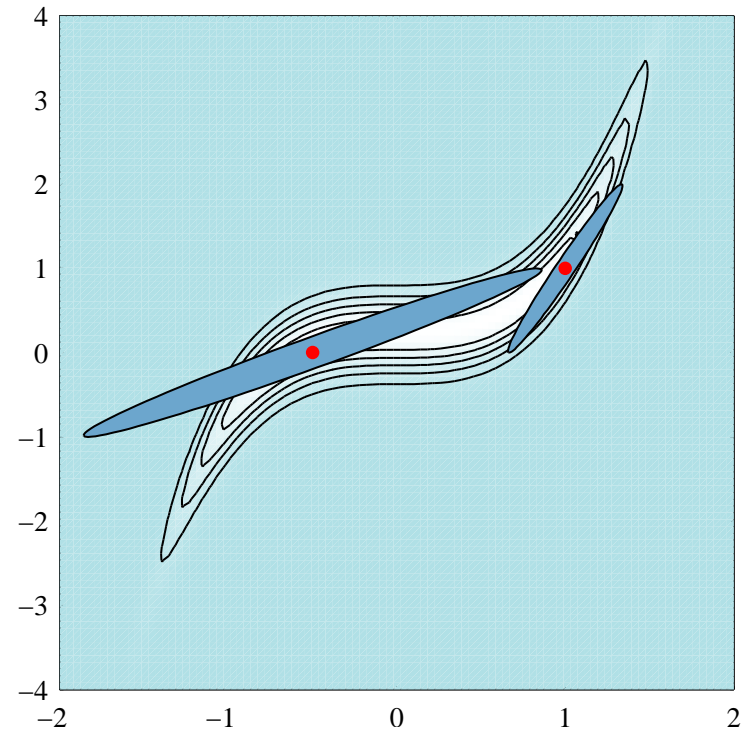


## ADAPTIVE METROPOLIS-HASTINGS (AM) ALGORITHM

*Adaptation:* As the sampling proceeds, the proposal distribution is updated to conform with the underlying density.



## ADAPTATION STRATEGIES

*Analytic adaptation:* Calculate a local Gaussian approximation of the probability density. Update every  $M$ th step.

May be computationally heavy

May be difficult to find analytically and requires numerical differentiation

*Sampling-based adaptation:* Use the already calculated sample history to determine the new proposal distribution.

## SAMPLING-BASED ADAPTATION

Design an algorithm along the following guidelines.

- Start a usual MH sampling at a point of your choice using a white noise proposal distribution.
- After possibly removing a burn-in sequence from the beginning, calculate the empirical covariance of the sample points obtained thus far.
- Use the empirical covariance to sample new points.
- Update the covariance every  $M$ th sample point.

## ADAPTED RANDOM WALK MH ALGORITHM

1. Initialize  $k = 0$ ,  $C_k = \gamma^2 I$ .

2. Generate a sample sequence of length  $M$ ,

$$x_{kM+1}, x_{kM+2}, \dots, x_{(k+1)M},$$

using the random walk proposal

$$x_{\text{prop}} = x_{\text{curr}} + w, \quad w \sim \mathcal{N}(0, C_k)$$

3. Update

$$C_k \rightarrow C_{k+1} = \text{cov}(x_1, x_2, \dots, x_{(k+1)M}) + \varepsilon I.$$

4. Increase  $k \rightarrow k + 1$  and continue from 2 until desired sample size is reached.

## QUESTIONS, ANSWERS

- Is this a Markov Chain method? The update depends on *all* of the sample history via the covariance matrix!

True! But one can show that little by little, the process forgets the past, and is *asymptotically* Markovian and therefore ergodic.

- What is that little  $\varepsilon I$  doing there?

It has two functions: Practical function is to avoid the possibility that all sample points become collinear. Theoretical function is to make sure that the ergodicity works.

- How do we update the covariance in practice?

Let's see...

## COVARIANCE UPDATING, STABLY

Divide the sample in blocks of length  $M$ :

$$\underbrace{x_1, x_2, \dots, x_M}_M, \underbrace{x_{M+1}, x_{M+2}, \dots, x_{2M}}_M, x_{2M+1}, \dots$$

Average and covariance over subsamples:

$$\hat{x}_k = \frac{1}{M} \sum_{j=(k-1)M+1}^{kM} x_j,$$

$$\hat{C}_k = \frac{1}{M} \sum_{j=(k-1)M+1}^{kM} (x_j - \hat{x}_k)(x_j - \hat{x}_k).$$

Denote the *cumulative* mean and covariance by

$$\bar{x}_{kM} = \frac{1}{kM} \sum_{j=1}^{kM} x_j,$$

$$\bar{C}_{kM} = \frac{1}{kM} \sum_{j=1}^{kM} (x_j - \bar{x}_{kM})(x_j - \bar{x}_{kM}).$$

Problem: *Find a numerically stable way of updating*

$$\bar{x}_{kM} \rightarrow \bar{x}_{(k+1)M},$$

$$\bar{C}_{kM} \rightarrow \bar{C}_{(k+1)M}.$$

## UPDATING THE MEAN

We have

$$\begin{aligned}\bar{x}_{(k+1)M} &= \frac{1}{(k+1)M} \sum_{j=1}^{(k+1)M} x_j \\ &= \frac{k}{k+1} \frac{1}{kM} \sum_{j=1}^{kM} x_j + \frac{1}{(k+1)M} \sum_{j=kM+1}^{(k+1)M} x_j \\ &= \frac{k}{k+1} \bar{x}_{kM} + \frac{1}{k+1} \hat{x}_{k+1}.\end{aligned}$$



## UPDATING THE COVARIANCE

We need some auxiliary results.

FACT 1: If

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j,$$

the covariance can be written as

$$\begin{aligned} C &= \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})(x_j - \bar{x})^T \\ &= \frac{1}{n} \sum_{j=1}^n x_j x_j^T - \underbrace{\frac{1}{n} \sum_{j=1}^n x_j \bar{x}^T}_{=\bar{x}} - \bar{x} \underbrace{\frac{1}{n} \sum_{j=1}^n x_j^T}_{=\bar{x}^T} + \bar{x} \bar{x}^T \\ &= \frac{1}{n} \sum_{j=1}^n x_j x_j^T - \bar{x} \bar{x}^T. \end{aligned}$$

FACT 2: For any  $\tilde{x}$ , the *non-centered covariance* is

$$\begin{aligned}
\tilde{C} &= \frac{1}{n} \sum_{j=1}^n (x_j - \tilde{x})(x_j - \tilde{x})^T \\
&= \frac{1}{n} \sum_{j=1}^n x_j x_j^T - \underbrace{\frac{1}{n} \sum_{j=1}^n x_j}_{=\bar{x}} \tilde{x}^T - \tilde{x} \underbrace{\frac{1}{n} \sum_{j=1}^n x_j^T}_{=\bar{x}^T} + \tilde{x} \tilde{x}^T \\
&= \frac{1}{n} \sum_{j=1}^n x_j x_j^T - \bar{x} \bar{x}^T + (\bar{x} \bar{x}^T - \bar{x} \tilde{x}^T - \tilde{x} \bar{x}^T + \tilde{x} \tilde{x}^T) \\
&= C + (\bar{x} - \tilde{x})(\bar{x} - \tilde{x})^T.
\end{aligned}$$

With these results, we have

$$\begin{aligned}
 C_{(k+1)M} &= \frac{1}{(k+1)M} \sum_{j=1}^{(k+1)M} (\mathbf{x}_j - \bar{\mathbf{x}}_{(k+1)M})(\mathbf{x}_j - \mathbf{x}_{(k+1)M})^T \\
 &= \frac{1}{(k+1)M} \left( \sum_{j=1}^{kM} + \sum_{kM+1}^{(k+1)M} \right) (\mathbf{x}_j - \bar{\mathbf{x}}_{(k+1)M})(\mathbf{x}_j - \mathbf{x}_{(k+1)M})^T.
 \end{aligned}$$

Both terms are, up to a multiplicative factor, non-centered covariances.

First term:

$$\begin{aligned}
& \frac{1}{(k+1)M} \sum_{j=1}^{kM} (x_j - \bar{x}_{(k+1)M})(x_j - \bar{x}_{(k+1)M})^T \\
&= \frac{k}{k+1} \frac{1}{kM} \sum_{j=1}^{kM} (x_j - \bar{x}_{kM})(x_j - \bar{x}_{kM})^T \\
&\quad + \frac{k}{k+1} (\bar{x}_{kM} - \bar{x}_{(k+1)M})(\bar{x}_{kM} - \bar{x}_{(k+1)M})^T \\
&= \frac{k}{k+1} C_{kN} + \frac{k}{k+1} (\bar{x}_{kM} - \bar{x}_{(k+1)M})(\bar{x}_{kM} - \bar{x}_{(k+1)M})^T.
\end{aligned}$$

Substituting the updating formula:

$$\begin{aligned}\bar{x}_{kM} - \bar{x}_{(k+1)M} &= \bar{x}_{kM} - \frac{k}{k+1}\bar{x}_{kM} - \frac{1}{k+1}\hat{x}_{k+1} \\ &= \frac{1}{k+1}(\bar{x}_{kM} - \hat{x}_{k+1}),\end{aligned}$$

so we have

$$\begin{aligned}&\frac{1}{(k+1)M} \sum_{j=1}^{kM} (x_j - \bar{x}_{(k+1)M})(x_j - \bar{x}_{(k+1)M})^T \\ &= \frac{k}{k+1}C_{kN} + \frac{k}{(k+1)^3}(\bar{x}_{kM} - \hat{x}_{k+1})(\bar{x}_{kM} - \hat{x}_{k+1})^T.\end{aligned}$$

Second term:

$$\begin{aligned}
& \frac{1}{(k+1)M} \sum_{j=kM+1}^{(k+1)M} (x_j - \bar{x}_{(k+1)M})(x_j - \bar{x}_{(k+1)M})^T \\
& \frac{1}{k+1} \frac{1}{M} \sum_{j=kM+1}^{(k+1)M} (x_j - \hat{x}_{k+1})(x_j - \hat{x}_{k+1})^T \\
& \quad + \frac{1}{k+1} (\hat{x}_{k+1} - \bar{x}_{(k+1)M})(\hat{x}_{k+1} - \bar{x}_{(k+1)M})^T \\
& \frac{1}{k+1} \hat{C}_{k+1} + \frac{1}{k+1} (\hat{x}_{k+1} - \bar{x}_{(k+1)M})(\hat{x}_{k+1} - \bar{x}_{(k+1)M})^T.
\end{aligned}$$

Again, substituting the recursion formula gives

$$\begin{aligned}\widehat{\boldsymbol{x}}_{k+1} - \bar{\boldsymbol{x}}_{(k+1)M} &= \widehat{\boldsymbol{x}}_{k+1} - \frac{k}{k+1}\bar{\boldsymbol{x}}_{kM} - \frac{1}{k+1}\widehat{\boldsymbol{x}}_{k+1} \\ &= \frac{k}{k+1}(\widehat{\boldsymbol{x}}_{k+1} - \bar{\boldsymbol{x}}_{kM}),\end{aligned}$$

and therefore

$$\begin{aligned}&\frac{1}{(k+1)M} \sum_{j=kM+1}^{(k+1)M} (\boldsymbol{x}_j - \bar{\boldsymbol{x}}_{(k+1)M})(\boldsymbol{x}_j - \bar{\boldsymbol{x}}_{(k+1)M})^T \\ &= \frac{1}{k+1}\widehat{\boldsymbol{C}}_{k+1} + \frac{k^2}{(k+1)^3}(\bar{\boldsymbol{x}}_{kM} - \widehat{\boldsymbol{x}}_{k+1})(\bar{\boldsymbol{x}}_{kM} - \widehat{\boldsymbol{x}}_{k+1})^T.\end{aligned}$$

Putting the pieces together gives

$$\begin{aligned}
C_{(k+1)M} &= \frac{k}{k+1} C_{kN} + \frac{k}{(k+1)^3} (\bar{x}_{kM} - \hat{x}_{k+1}) (\bar{x}_{kM} - \hat{x}_{k+1})^T \\
&+ \frac{1}{k+1} \hat{C}_{k+1} + \frac{k^2}{(k+1)^3} (\bar{x}_{kM} - \hat{x}_{k+1}) (\bar{x}_{kM} - \hat{x}_{k+1})^T \\
&= \frac{k}{k+1} C_{kN} + \frac{1}{k+1} \hat{C}_{k+1} + \frac{k}{(k+1)^2} (\bar{x}_{kM} - \hat{x}_{k+1}) (\bar{x}_{kM} - \hat{x}_{k+1})^T,
\end{aligned}$$

which is the desired updating formula.



EXAMPLE: PROOF OF CONCEPT

Sampling a Gaussian in  $\mathbb{R}^2$ .

$$\pi(x) \propto \exp\left(-\frac{1}{2}(x-b)^T\Gamma^{-1}(x-b)\right),$$

where

$$b = \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \quad \Gamma = UDU^T,$$

$$U = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}, \quad \theta = \frac{\pi}{3},$$

$$D = \text{diag}(1, 0.1).$$

## THE PROGRAM: ADAPTATION VS. NON-ADAPTATION

```
% Defining the underlying distribution: a 2D Gaussian
```

```
th = pi/3;  
U = [cos(th), -sin(th); sin(th), cos(th)];  
d = [1, 0.1];  
D = diag(d);  
Gamma = U*D*U';  
b = [2; 2];  
invGamma = inv(Gamma);
```

```
% Initializing
```

```
x0 = [3;1];    % Initial sampling point  
step = 0.02;  % Initial step: no prior tuning for MH  
tiny = 1e-6;  
nsample = 150000;  
M = 100;      % Adaptation period
```

Observe: the step size is way too small for non-adaptive MH. The point here is to demonstrate that the adaptive method requires no tuning, i.e., you can start with sub-optimal proposal.

```
% Sampling without adaptation

SampleNA = zeros(2,nsample);
SampleNA(:,1) = x0;
x = x0;
lpdf = -0.5*(x-b)'*invGamma*(x-b);
accrate = 0;
```

```

for j = 2:nsample
    % Draw the proposal
    xprop = x + step*randn(2,1);
    lpdfprop = -0.5*(xprop-b)'invGamma*(xprop-b);
    % Check for acceptance
    if lpdfprop - lpdf >log(rand)
        %accept
        x = xprop;
        lpdf = lpdfprop;
        accrate = accrate + 1;
    end
    SampleNA(:,j) = x;
end
rel_accrate = 100*accrate/nsample

```

## SAMPLING WITH ADAPTATION

Updating

$$x_{\text{prop}} = x_{\text{curr}} + s, \quad s \sim \mathcal{N}(0, C),$$

that is,

$$\pi(s) \propto \exp\left(-\frac{1}{2}s^{\text{T}}C^{-1}s\right).$$

Write the Cholesky decomposition,

$$C = R^{\text{T}}R,$$

so

$$C^{-1} = R^{-1}R^{-\text{T}}.$$

This means that

$$\pi(s) \propto \exp\left(-\frac{1}{2}\|R^{-\text{T}}s\|^2\right),$$

or that

$$R^{-\text{T}}s = w \sim \mathcal{N}(0, I).$$

Hence, the updating procedure is

$$x_{\text{prop}} = x_{\text{curr}} + R^{\text{T}}w, \quad w \sim \mathcal{N}(0, I).$$

```
% Sampling with adaptation

SampleA = zeros(2,nsample);
SampleA(:,1) = x0;
x = x0;
lpdf = -0.5*(x-b)'invGamma*(x-b);
C = step^2*eye(2);
mean = zeros(2,1);
R = step*eye(2);
accrate = 0;
tempSample = [x];
k = 0;
```



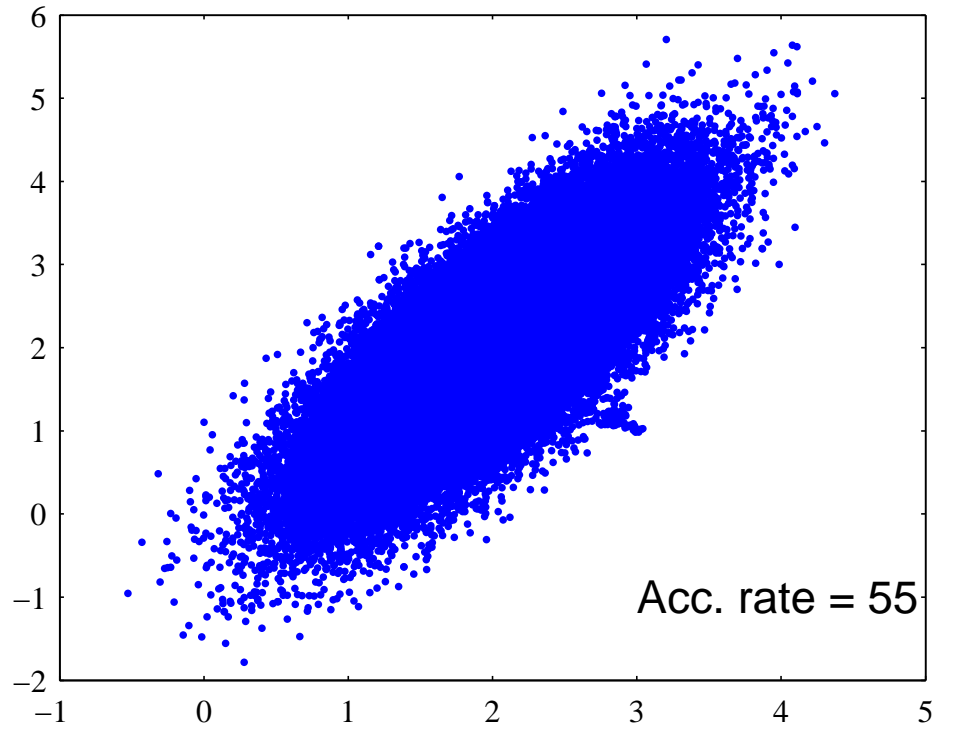
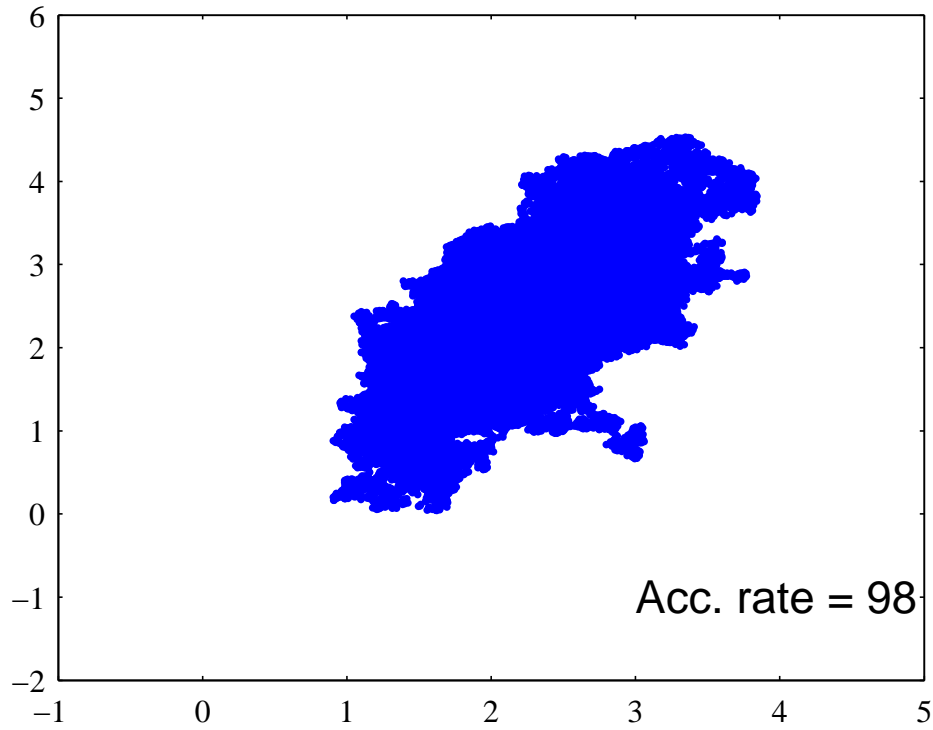
```
for j =2:nsample
    % Draw the proposal
    xprop = x + R'*randn(2,1);
    lpdfprop = -0.5*(xprop-b)'invGamma*(xprop-b);
    % Check for acceptance
    if lpdfprop - lpdf >log(rand)
        %accept
        x = xprop;
        lpdf = lpdfprop;
        accrate = accrate + 1;
    end
    SampleA(:,j) = x;
    tempSample = [tempSample x];
end
```

```

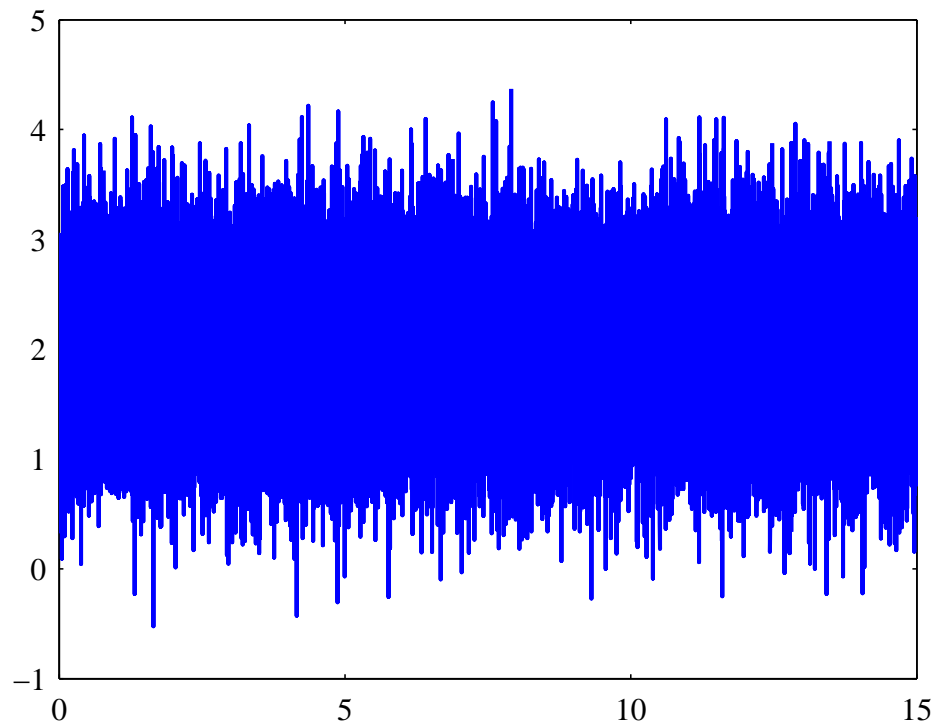
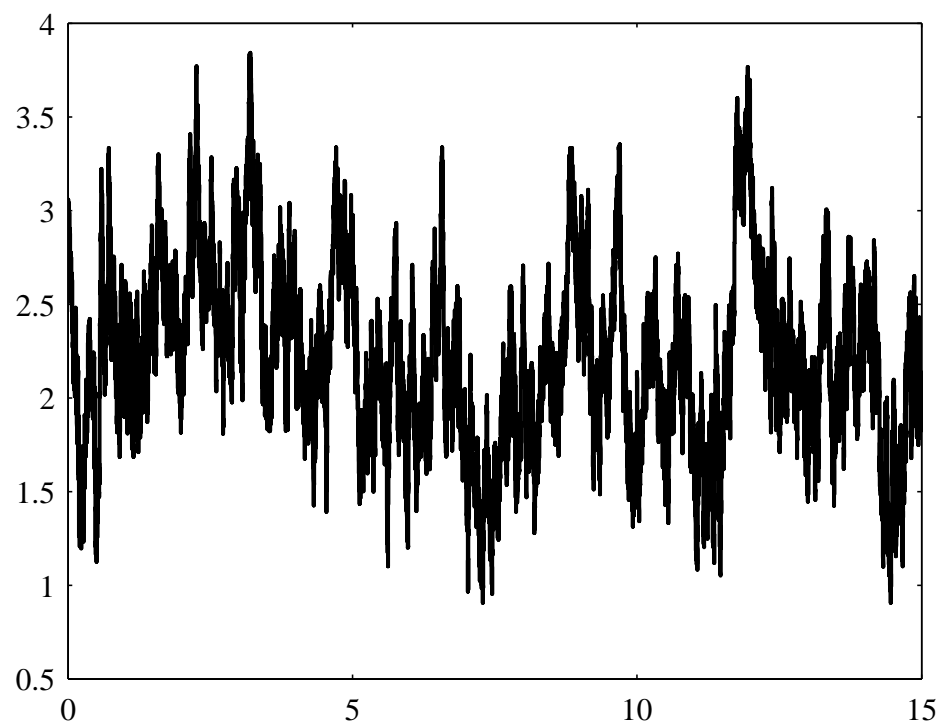
if mod(j,M) == 0
    % Update the proposal distribution
    xk = 1/M*sum(tempSample')';
    aux = tempSample - xk*ones(1,M);
    Ck = 1/M*aux*aux';
    mean = k/(k+1)*mean + 1/(k+1)*xk;
    C = k/(k+1)*C + 1/(k+1)*Ck + k/(k+1)^2*(mean-xk)*(mean-xk)';
    R = chol(C);
    k = k+1;
    tempSample = [];
end
end rel_accrateA = 100*accrate/nsample

```

# SCATTER PLOTS



# SAMPLE HISTORIES



## DIAGNOSTICS

The  $p\%$  probability region:

$$\Gamma = UDU^T \Rightarrow \Gamma^{-1} = UD^{-1}U^T,$$

$$\begin{aligned}\pi(x) &\propto \exp\left(-\frac{1}{2}(x-b)^T\Gamma^{-1}(x-b)\right) \\ &= \exp\left(-\frac{1}{2}\|D^{-1/2}U^T(x-b)\|^2\right),\end{aligned}$$

so

$$W = D^{-1/2}U^T(X-b) \sim \mathcal{N}(0, I).$$

The  $p\%$  probability region for  $W$  is a disc  $D_\alpha$  of radius  $\alpha$ :

$$\begin{aligned} \mathbb{P}\{W \in D_\alpha\} &= \frac{1}{2\pi} \int_{D_\alpha} \exp\left(-\frac{1}{2}\|w\|^2\right) dw \\ &= \frac{1}{2\pi} \int_0^\alpha \int_0^{2\pi} e^{-r^2/2} d\theta r dr \\ &= 1 - e^{-\alpha^2/2} = \frac{p}{100}, \end{aligned}$$

which is equivalent to

$$\alpha = \sqrt{2 \log \frac{100}{100 - p}}.$$

Given a sample  $\{x_1, x_2, \dots, x_N\}$ , we

- calculate the sample  $\{x_1, x_2, \dots, x_N\}$ ,

$$w_j = D^{-1/2}U^T(x_j - b),$$

- calculate the relative amount of these sample points are within the disc  $D_\alpha$ ,

$$r_p(N) = \frac{100}{N} \#\{w_j \mid \|w_j\| < \alpha\}.$$

When  $N$  grows, we should have

$$r_p(N) \rightarrow p.$$

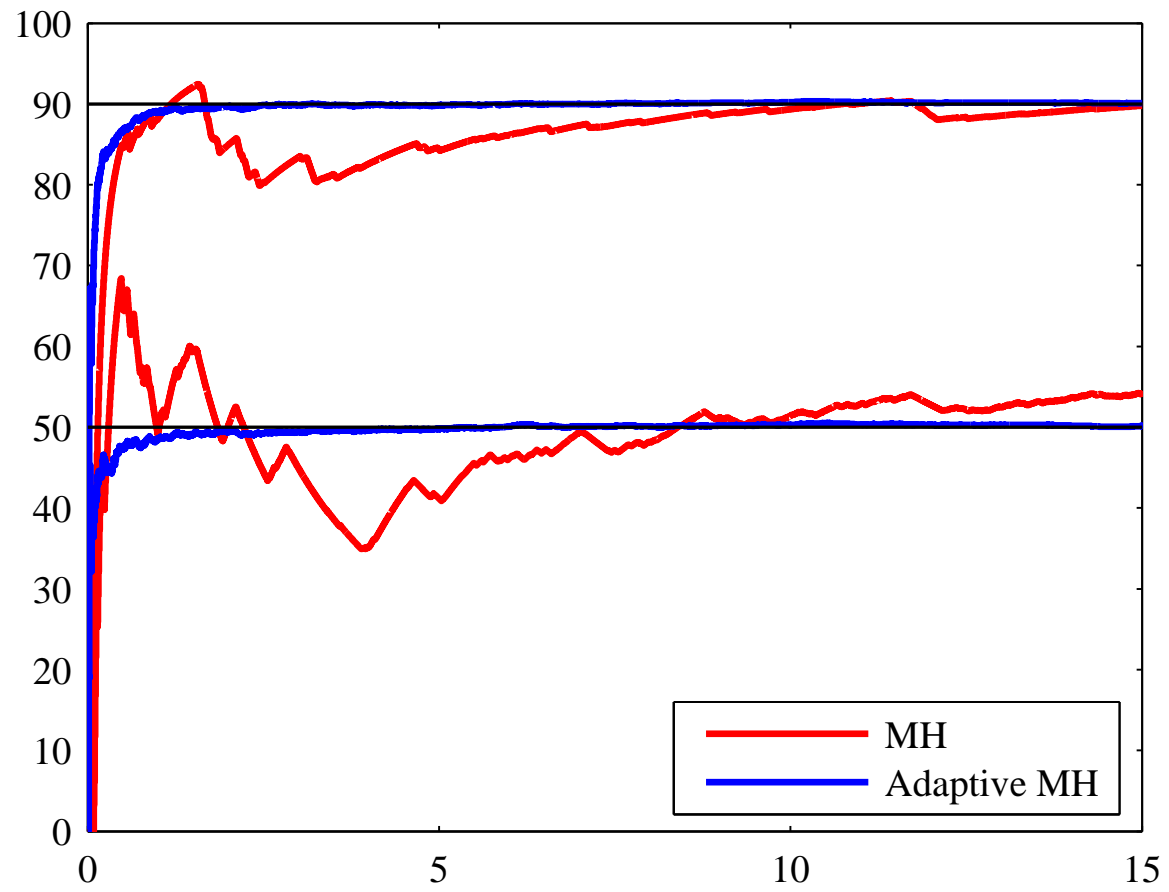
## PROGRAM

```
% Number of points within a p percent ellipse

W = diag(1 ./sqrt(d))*U'*(SampleNA-b*ones(1,nsample));
normW = sqrt(sum(W.^2));
p = [90,50];

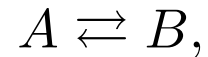
for j = 1:2
    alpha = sqrt(2*log(100/(100-p(j))));
    xinside = (normW<alpha);
    reln = 100*cumsum(xinside)./[1:nsample];
end
```





## EXAMPLE: INVERSE PROBLEMS IN CHEMICAL ENGINEERING

Recall the reversible chemical reactions



with reaction rates  $k_1$  and  $k_2$ , respectively.

Concentrations  $C_A$  and  $C_B$  satisfy

$$\frac{dC_A}{dt} = -k_1 C_A + k_2 C_B$$

$$\frac{dC_B}{dt} = k_1 C_A - k_2 C_B,$$

with initial data

$$C_A(0) = C_{A,0}, \quad C_B(0) = C_{b,0}.$$

## INVERSE PROBLEM

Assume that we know the initial concentrations.

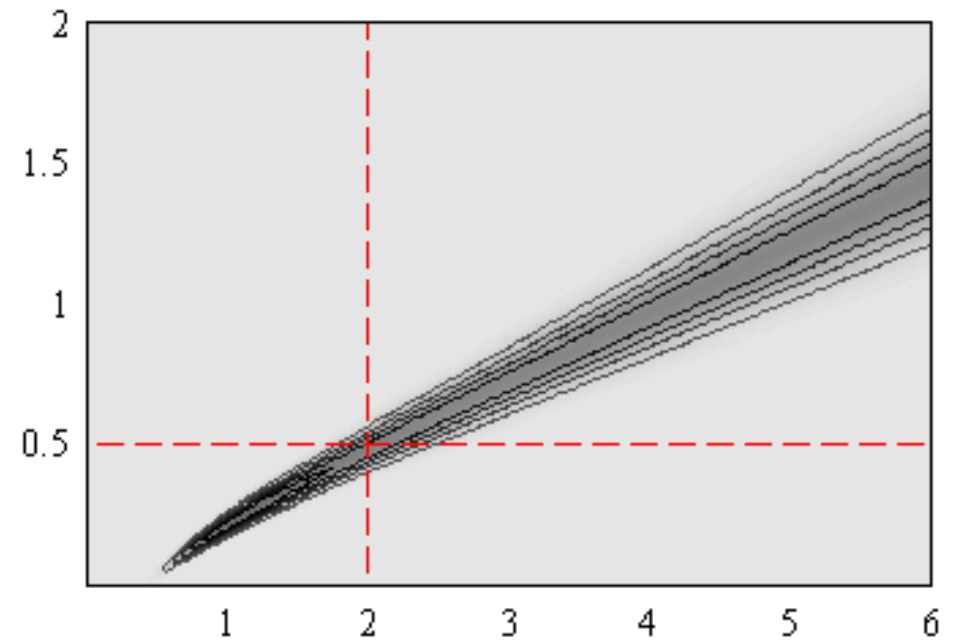
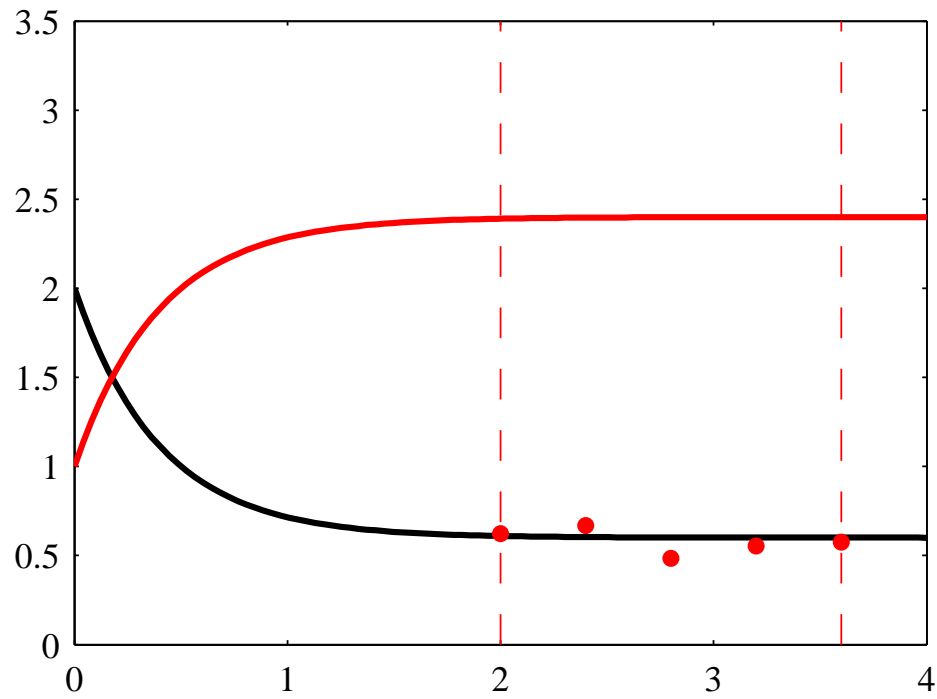
Data: For  $0 < t_1 < t_2 \cdots < t_n$ , measure  $C_A(t_j)$ ,  $1 \leq j \leq n$ .

Estimate  $k_1$  and  $k_2$ .

Noisy observations:

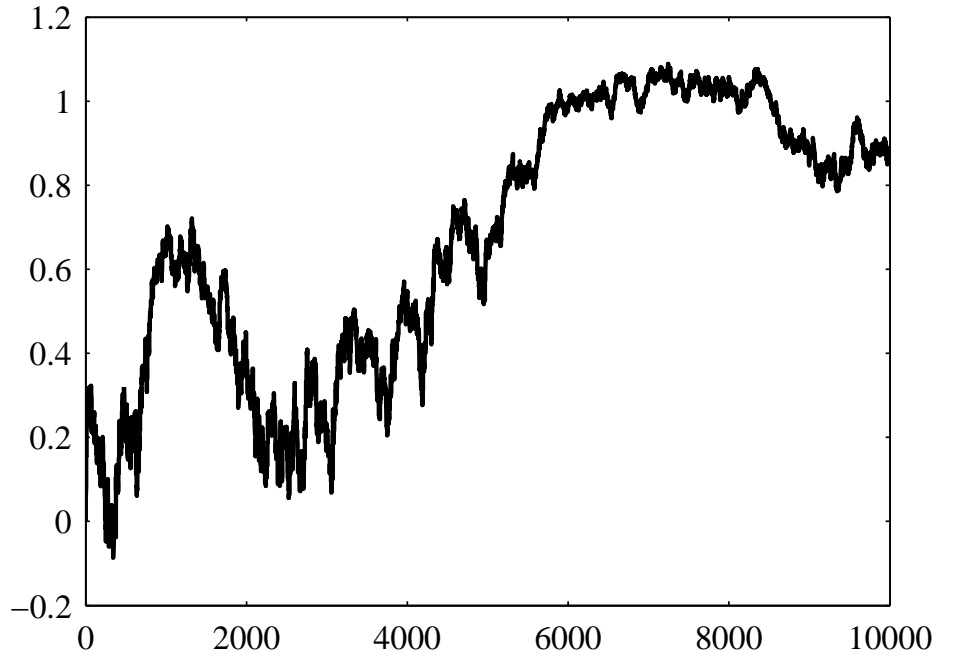
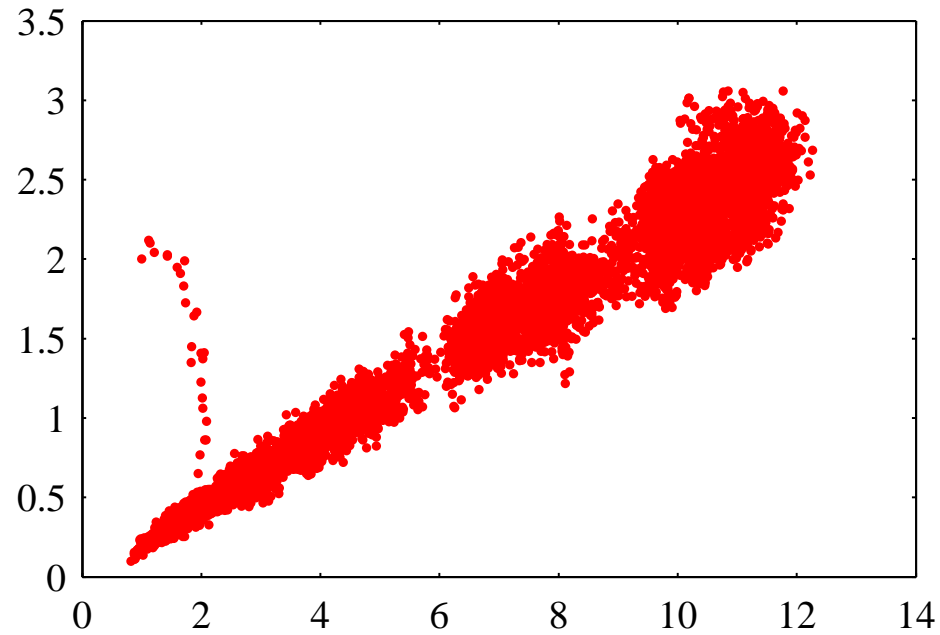
$$b_j = C_A(t_j) + e_j, \quad e_j \sim \mathcal{N}(0, \sigma^2).$$

# STEADY STATE DATA



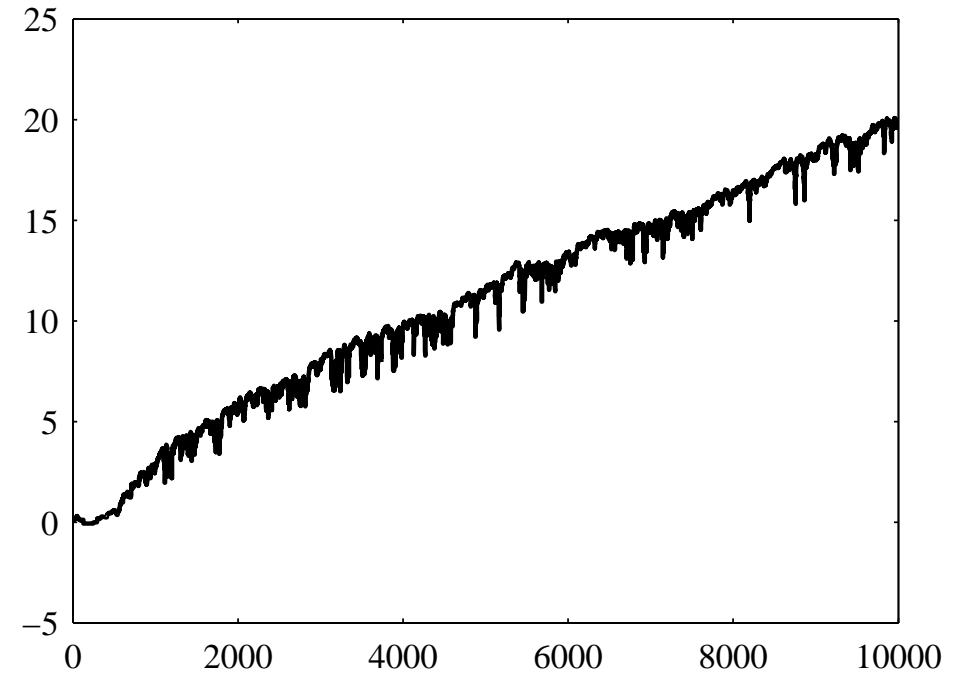
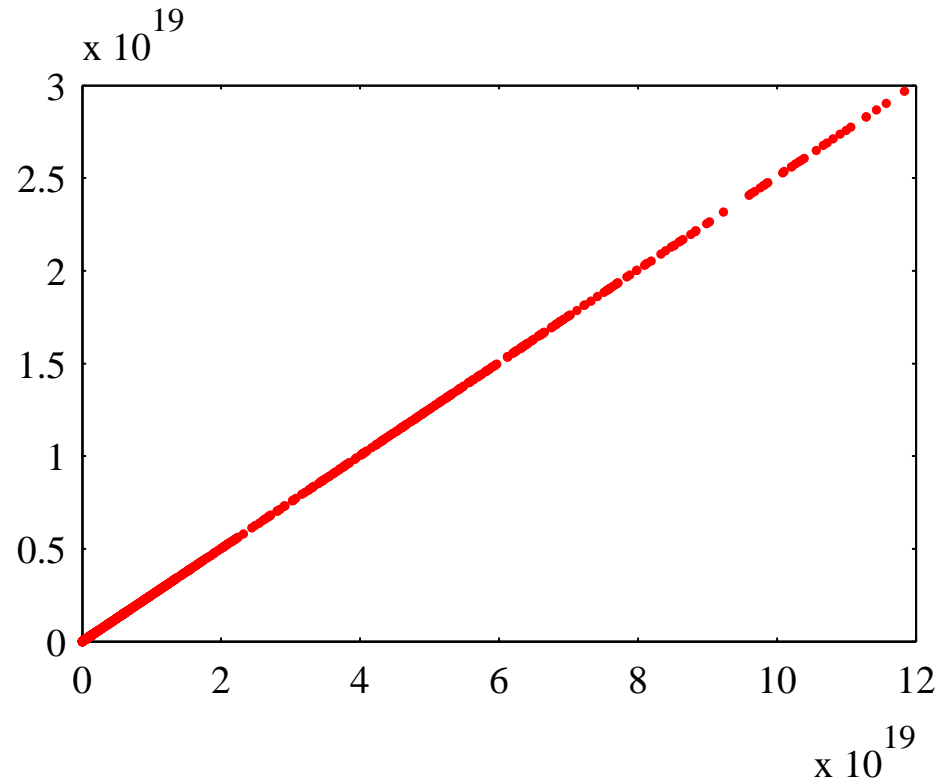
$$C_{A,0} = 2, \quad C_{B,0} = 1, \quad k_1 = 2, \quad k_2 = 0.5, \quad \sigma = 0.2$$

# SAMPLING WITH NON-ADAPTIVE MH



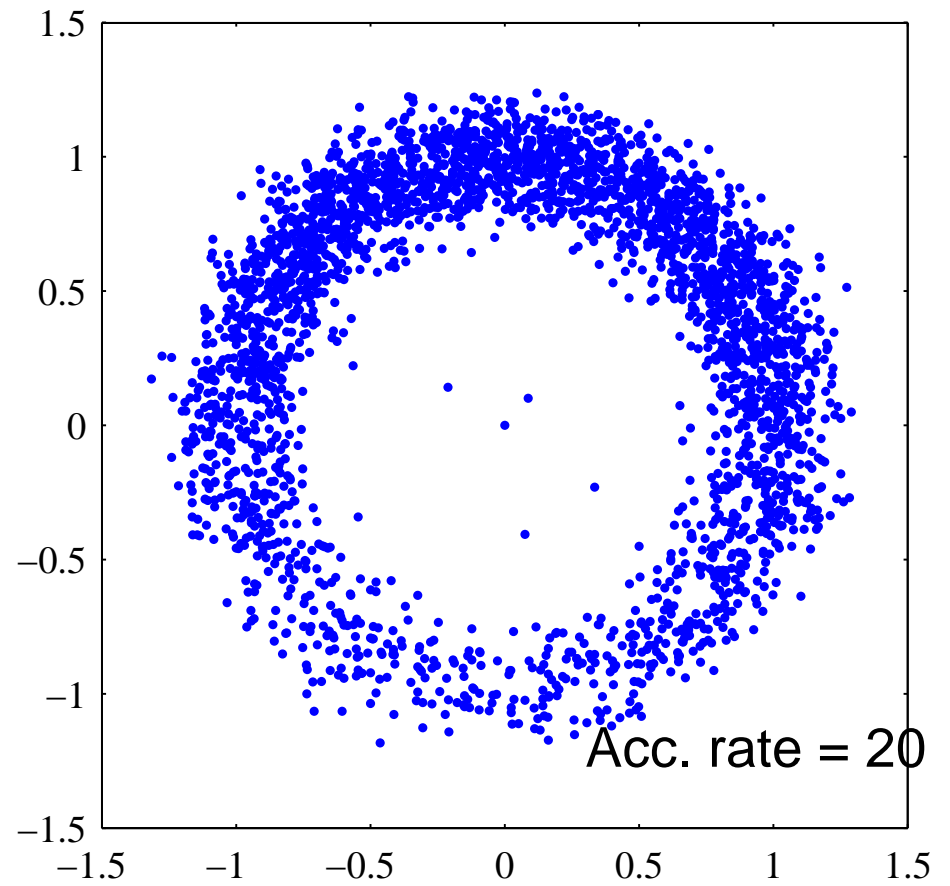
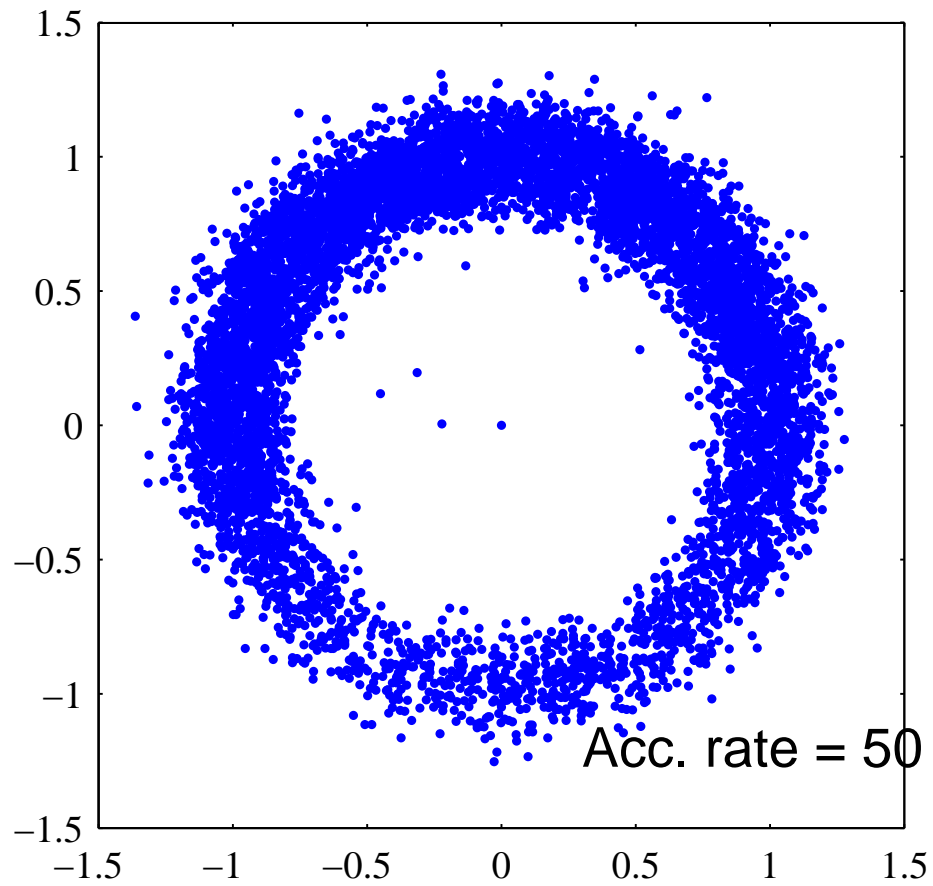
Initial point (1, 2)

# SAMPLING WITH ADAPTIVE MH

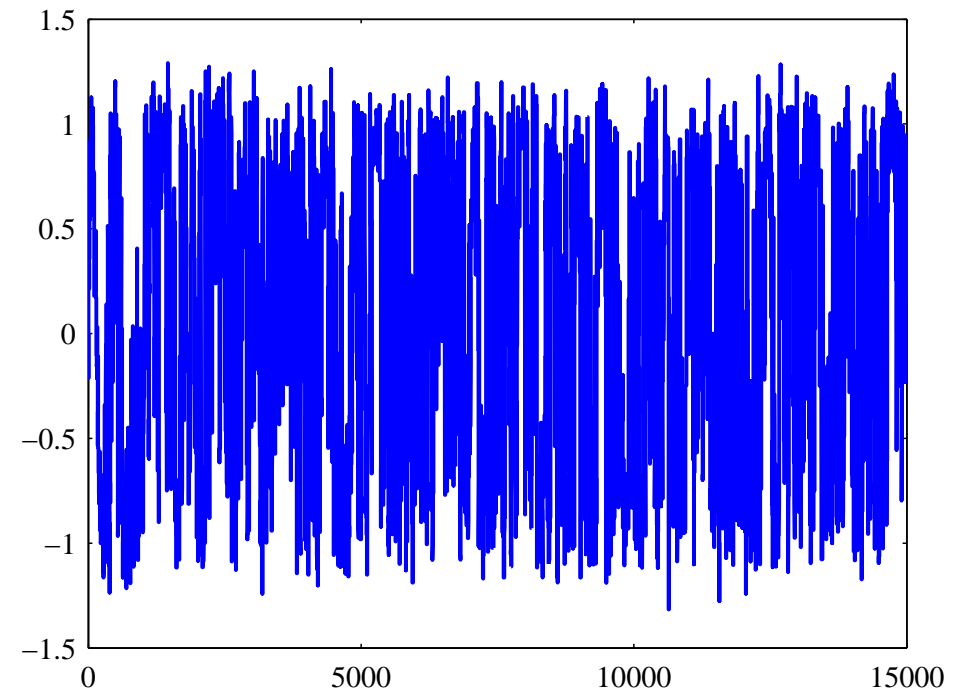
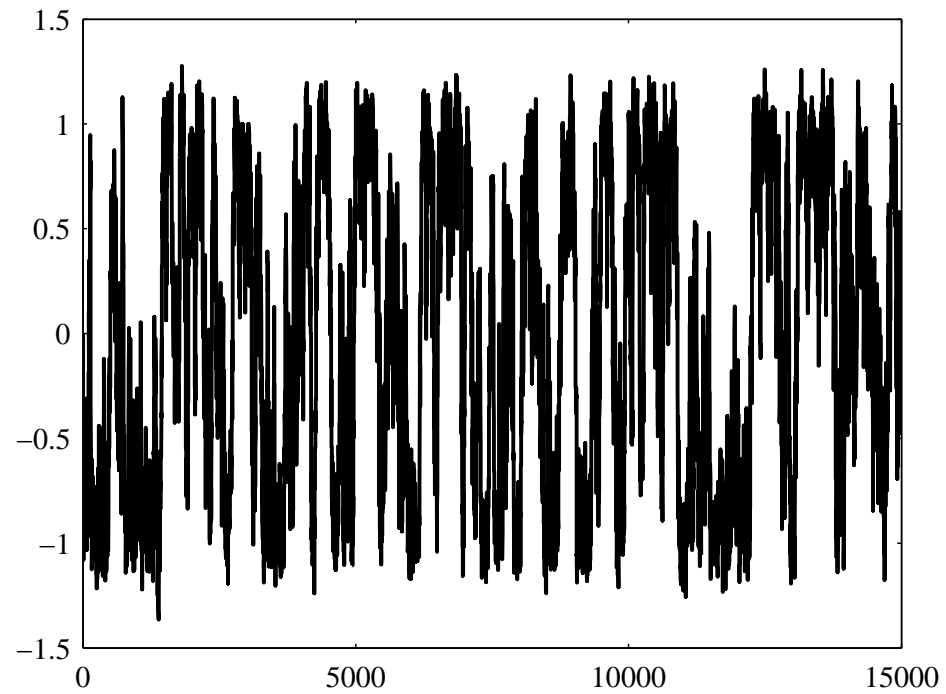


Adaptation after every 100 sample points

# EXAMPLE: HORSESHOE DISTRIBUTION



Non-adaptive (left) vs. adaptive (right).



Non-adaptive (left) vs. adaptive (right).



## OBSERVATIONS

In this example, the adaptation, as defined here, is not of great help:

The distribution is almost circular, so the asymptotic covariance is almost an identity, and we end up drawing essentially from a white noise density.

The only advantage is that the step length need no tuning.