# UiO : Matematisk institutt

Det matematisk-naturvitenskapelige fakultet

## STK-4051/9051  Computational Statistics  Spring 2021
## Markov Chain Monte Carlo

Instructor: Odd Kolbjørnsen, oddkol@math.uio.no

# Last time variance reduction

- Beating: $\dfrac{\mathrm{Var}\{h(X)\}}{n}$

- Antithetic sampling

  - Random numbers that have negative correlation

Define $\hat{\mu}_{AS} = \frac{1}{2}(\hat{\mu}_1 + \hat{\mu}_2)$

$\mathrm{var}[\hat{\mu}_{AS}] = \frac{1}{4}(\mathrm{var}[\hat{\mu}_1] + \mathrm{var}[\hat{\mu}_2]) + \frac{1}{2}\mathrm{cov}[\hat{\mu}_1, \hat{\mu}_2]$

$\quad\quad\quad = \frac{(1+\rho)\sigma^2}{2n}$

- Exercise: Common random numbers

  - Creating a paired test rather than a two sample distribution (when appropriate)

- Importance sampling

  - Normalized weights vs un-normalized

- Control variates

$\hat{\mu}_{CV} = \hat{\mu}_{MC} + \lambda(\hat{\theta}_{MC} - \theta)$

$\mathrm{var}[\hat{\mu}_{CV}] = \mathrm{var}[\hat{\mu}_{MC}] + \lambda^2 \mathrm{var}[\hat{\theta}_{MC}] + 2\lambda\mathrm{cov}[\hat{\mu}_{MC}, \hat{\theta}_{MC}]$

  - We know something about the distribution

$\lambda = -\dfrac{\mathrm{cov}[\hat{\mu}_{MC}, \hat{\theta}_{MC}]}{\mathrm{var}[\hat{\theta}_{MC}]}$

- Rao-Blacwellization

  - We know something about a conditional distribution
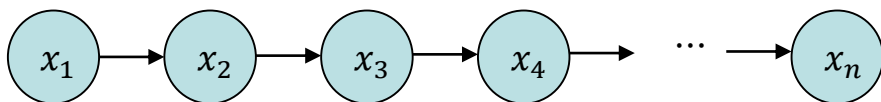
  - Particular useful with hyper parameters

$\mathrm{var}[h(\mathbf{X}_i)] = E[\mathrm{var}[h(\mathbf{X}_i)|\mathbf{X}_2]] + \mathrm{var}[E[h(\mathbf{X})|\mathbf{X}_2]] \geq \mathrm{var}[E[h(\mathbf{X})|\mathbf{X}_2]]$

# Today

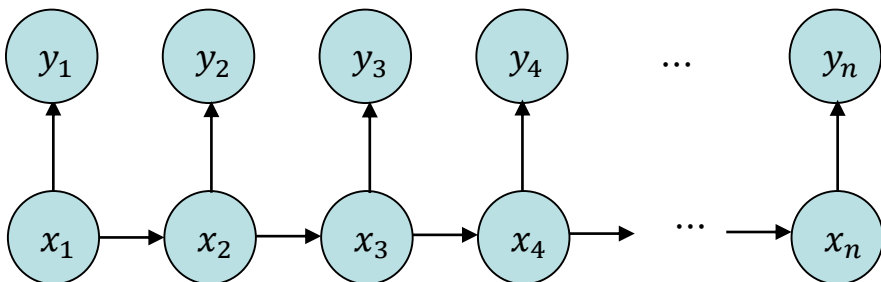- <span style="color:red">Exact</span> methods
  - Inversion/transformation methods
  - Rejection sampling
- <span style="color:red">Approximate</span> methods
  - Sampling importance resampling
  - Sequential Monte Carlo
  - Markov chain Monte Carlo (Chapter 7 and 8)
- <span style="color:red">Variance reduction</span> methods
  - Importance sampling
  - Antithetic sampling
  - Control variates
  - Rao-blackwellization
  - Common random numbers

# Graphing the probability distribution

The way I use it is to highlight the dependency
structure in a statistical model model, i.e. the joint didtribution

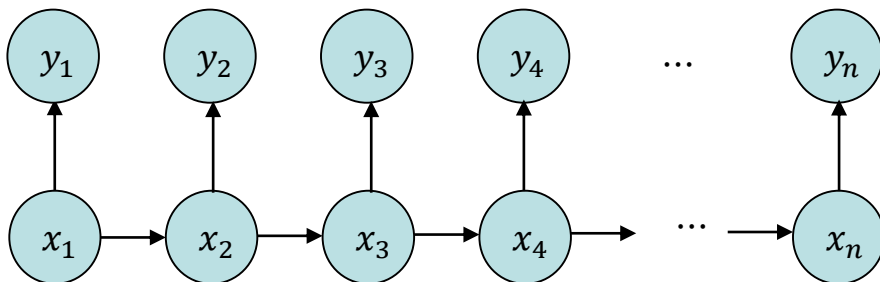$$f(\boldsymbol{x}) = f(x_1)f(x_2|x_1)f(x_3|x_2) \cdots f(x_n|x_{n-1})$$

$$f(\boldsymbol{x}, \boldsymbol{y}) = f(x_1)f(y_1|x_1)f(x_2|x_1)f(y_2|x_2) \cdots f(x_n|x_{n-1})f(y_n|x_n)$$

# In a hidden Markov model

- The way I understand this is that the y's are observed and that the x's are hidden (unknown)

- What do we mean by saying that the $x_i$, shadows for $y_i$?
  - We mean this in the sense of conditional distributions

  When $x_1$ shadows for $y_1$ (wrt $x_2$ ) we have: $f(x_2|x_1, y_1) = f(x_2|x_1)$



In the graph the only way information from $y_1$ may get to $x_2$ is through its influence on $x_1$, thus if we know the value of $x_1$ , then there is no additional effect of $y_1$ on $x_2$

6

# **About SCM**

- In the lecture about SMC recap slide 17. We want to compute $p(\boldsymbol{y}_s | \boldsymbol{y}_{1:(s-1)}, \theta)$.
  We do this by integrating out $x_s$, but these variables are hidden, i.e. Data we do not have. How can we do that?

- We do **not** know the distribution $p(\boldsymbol{y}|\theta)$

- We **know** joint distribution of $p(\boldsymbol{x}, \boldsymbol{y}|\theta)$

  – So we know something about the x's (but not the value)

- Since we have not observed the $\boldsymbol{x}$'s we need to get rid of it, i.e. integrating it out.

$$p(\boldsymbol{y}|\theta) = \int p(\boldsymbol{x}, \boldsymbol{y}|\theta) d\boldsymbol{x}$$

# **What is a "sufficient statistic"?**

- A statistic is a function of the data, i.e. $S(\boldsymbol{X})$.

  – e.g.   $\frac{1}{n}\sum_{i=1}^{n} X_i$

- Given a model with parameters, a set of statistics is said to be sufficient  for a parameter $\theta$ if the distribution of data $\boldsymbol{X}$ conditioned to the statistics $S(\boldsymbol{X})$ do not depend on the parameter, $\theta$.

  – e.g. $E(X_j) = \mu$  vs  $E\left(X_j \mid \frac{1}{n}\sum_{i=1}^{n} X_i\right) = \frac{1}{n}\sum_{i=1}^{n} X_i$

  For the normal distibution   $E\left(X_j \mid \frac{1}{n}\sum_{i=1}^{n} X_i = a\right) = a$

# **Sufficient statistic in stk 4051**

- EM in exponential family
  - Exponential family is linked to sufficient statistics
  - When we precompute properties the sufficient statistics comes into play

- SCM parameter estimation (Bayesian)
  - The sufficient statistics lets us decouple the parameter and the data
  - $p(\boldsymbol{x}, \boldsymbol{s}, \boldsymbol{\theta}) = p(\boldsymbol{x}|\boldsymbol{s})p(\boldsymbol{s}|\boldsymbol{\theta})p(\boldsymbol{\theta})$

# EM in exponential family

- The Exponential family:

$$f_y(\mathbf{y}|\boldsymbol{\theta}) = c_1(\mathbf{y})c_2(\boldsymbol{\theta})\exp\{\boldsymbol{\theta}^T\mathbf{s}(\mathbf{y})\}$$

- $\mathbf{s}(\mathbf{y})$ is a sufficient statistic:

$$
\begin{aligned}
f_s(\mathbf{s}|\boldsymbol{\theta}) &= \int_{\mathbf{y}:\mathbf{s}(\mathbf{y})=\mathbf{s}} f_y(\mathbf{y}|\boldsymbol{\theta})d\mathbf{y} \\
&= \int_{\mathbf{y}:\mathbf{s}(\mathbf{y})=\mathbf{s}} c_1(\mathbf{y})c_2(\boldsymbol{\theta})\exp\{\boldsymbol{\theta}^T\mathbf{s}(\mathbf{y})\}d\mathbf{y} \\
&= c_2(\boldsymbol{\theta})\exp\{\boldsymbol{\theta}^T\mathbf{s}\}\int_{\mathbf{y}:\mathbf{s}(\mathbf{y})=\mathbf{s}} c_1(\mathbf{y})d\mathbf{y} \\
&= c_2(\boldsymbol{\theta})\exp\{\boldsymbol{\theta}^T\mathbf{s}\}g(\mathbf{s})
\end{aligned}
$$

$$f(\mathbf{y}|\mathbf{s};\boldsymbol{\theta}) = \frac{f_y(\mathbf{y}|\boldsymbol{\theta})}{f_s(\mathbf{s}|\boldsymbol{\theta})} = \frac{c_1(\mathbf{y})c_2(\boldsymbol{\theta})\exp\{\boldsymbol{\theta}^T\mathbf{s}\}}{c_2(\boldsymbol{\theta})\exp\{\boldsymbol{\theta}^T\mathbf{s}\}g(\mathbf{s})} = \frac{c_1(\mathbf{y})}{g(\mathbf{s})}$$

The distribution of the data **y** given
the sufficient statistic **s** and parameter $\theta$

does not depend on
the parameter $\theta$

10

# Sufficient statistics for parameter estimation in SCM

- The distribution of data $X$ conditioned to the statistics $S(X)$ do not depend on the parameter, $\theta$.

- $p(x, s, \theta) = p(x|s)p(s|\theta)p(\theta)$



- Which gives:
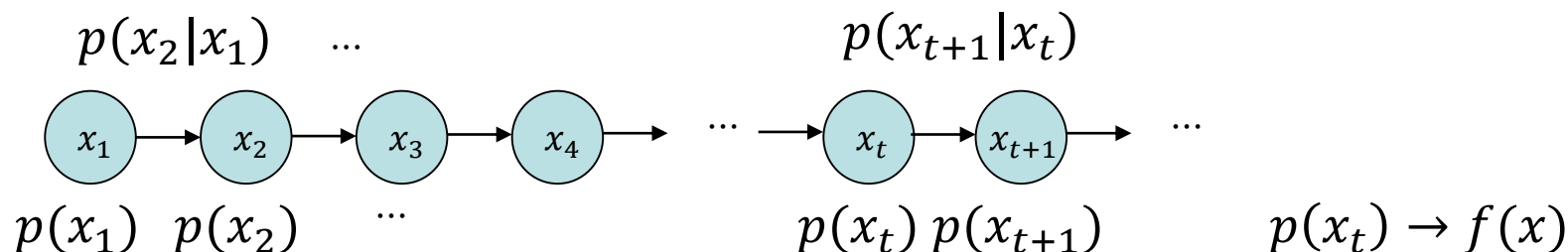    $$p(\theta|x, s) = p(\theta|s)$$

- We also have
    $$p(s|x, \theta) = p(s|x) = \delta(s = S(x))$$

# Questions

- Then we continiue on the topic of today

# Markov chain Monte Carlo

- Previously we computed weights to correct the distribution (or used rejection sampling)

- Now we will create a sequence of samples which will converge to samples from the correct distribution

$p(x_2|x_1)$ ...          $p(x_{t+1}|x_t)$

$$x_1 \rightarrow x_2 \rightarrow x_3 \rightarrow x_4 \rightarrow \quad \cdots \rightarrow x_t \rightarrow x_{t+1} \rightarrow \quad \cdots$$

$p(x_1)$ $p(x_2)$ ...          $p(x_t)\, p(x_{t+1})$          $p(x_t) \rightarrow f(x)$

# Markov chain Monte Carlo  (McMC)

- Assume now simulating from $f(\mathbf{X})$ is difficult directly
  - $f(\cdot)$ complicated
  - $\mathbf{X}$ high-dimensional
- Markov chain Monte Carlo:
  - Generates $\{\mathbf{X}^{(t)}\}$ sequentially
  - Markov structure: $\mathbf{X}^{(t)} \sim P(\cdot | \mathbf{X}^{(t-1)})$
- Aim now:
  - The distribution of $\mathbf{X}^{(t)}$ converges to $f(\cdot)$ as $t$ increases
  - $\hat{\mu}_{MCMC} = N^{-1} \sum_{t=1}^{N} h(\mathbf{X}^{(t)})$ converges towards $\mu = E^{f}[h(\mathbf{X})]$ as $t$ increases

Why?
We had problems with weight decay and degeneracy in the
direct approach now we can iterate to improve results

# Markov chain theory – discrete case

- Assume $\{X^{(t)}\}$ is a Markov chain where $X^{(t)}$ is a discrete random variable

  $$\Pr(X^{(t)} = y | X^{(t-1)} = x) = P(y|x)$$

  giving the transition probabilities
- Assume the chain is
  - irreducible: It is possible to move from any **x** to any **y** in a finite number of steps
  - reccurent: The chain will visit any state infinitely often.
  - aperiodic: Does not go in cycles
- Then there exists a unique distribution $f(x)$ such that

  $$\lim_{t\to\infty} \Pr(X^{(t)} = y | X^{(0)} = x) = f(y)$$

  $$\hat{\mu}_{MCMC} \to \mu = E^f[X]$$

  Limit distribution

- How to find $f(\cdot)$ (the stationary distribution): Solve

  $$f(y) = \sum_{x} f(x)P(y|x)$$

  Stationary distribution
  (fix point)

- Our situation: We have $f(y)$, want to find $P(y|x)$
  - Note: Many possible $P(y|x)$!

# Discrete Transition probability

$p(x_2|x_1)$ ...

$p(x_{t+1}|x_t)$

$$x_1 \rightarrow x_2 \rightarrow x_3 \rightarrow x_4 \rightarrow \quad \cdots \rightarrow x_t \rightarrow x_{t+1} \rightarrow \quad \cdots$$

$p(x_1)\ p(x_2)$ ...

$p(x_t)\ p(x_{t+1})$

$p(x_t) \rightarrow f(x)$

- Need initial distribution $p(x_1)$, say we have 4 possible classes
- and transition probability $p(x_t|x_{t-1})$, we need a transition to each state

|  |  |  |  |  |  | $x_2$ |
|  |  |  |  |  |  |  |
| $x_1$ | $p(x_1)$ |  | 1 | 2 | 3 | 4 |
| 1 | 0.4 | $p(x_2|x_1 = 1)$ | 0.80 | 0.10 | 0.00 | 0.10 |
| 2 | 0.1 | $p(x_2|x_1 = 2)$ | 0.05 | 0.90 | 0.05 | 0.00 |
| 3 | 0.1 | $p(x_2|x_1 = 3)$ | 0.00 | 0.05 | 0.90 | 0.05 |
| 4 | 0.4 | $p(x_2|x_1 = 4)$ | 0.10 | 0.00 | 0.10 | 0.80 |

$p_0 =$

$P =$

16

$p_0$

$p_0 \quad p_0^2 P$

$p(x_{20})$

| | |
|---|---|
| | 0.17 |
| | 0.33 |
| | 0.33 |
| | 0.17 |

$p_0 \quad p_0 P \qquad \ldots \qquad p_0 P^7 \qquad \ldots \qquad p_0 P^{19}$
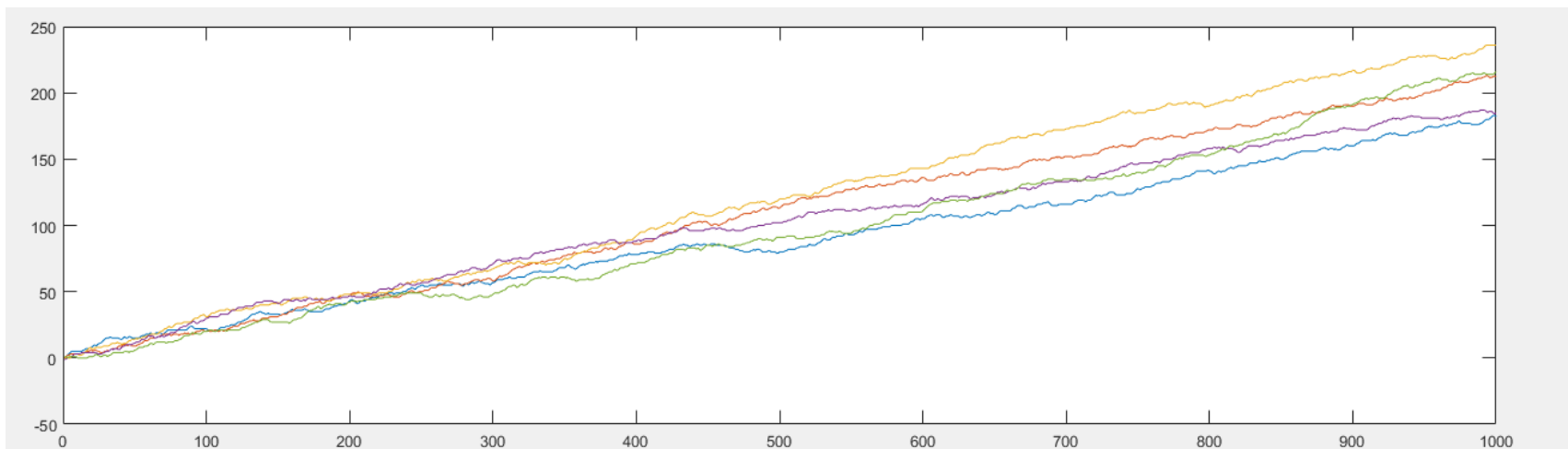
# Irreducible/ aperiodic:
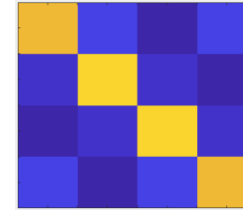
# **Recurrent (OK if finite and irreducible)**

- ## Problem if countable many discrete classes

$$P(x_t|x_{t-1}) = \begin{cases} 0.6 & x = x_{t-1} \\ 0.3 & x = x_{t-1} + 1 \\ 0.1 & x = x_{t-1} - 1 \end{cases}$$
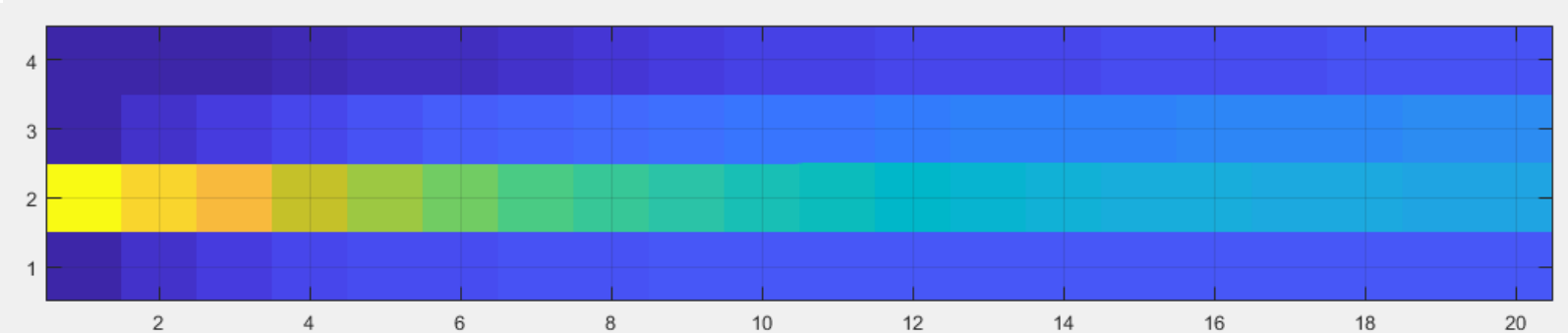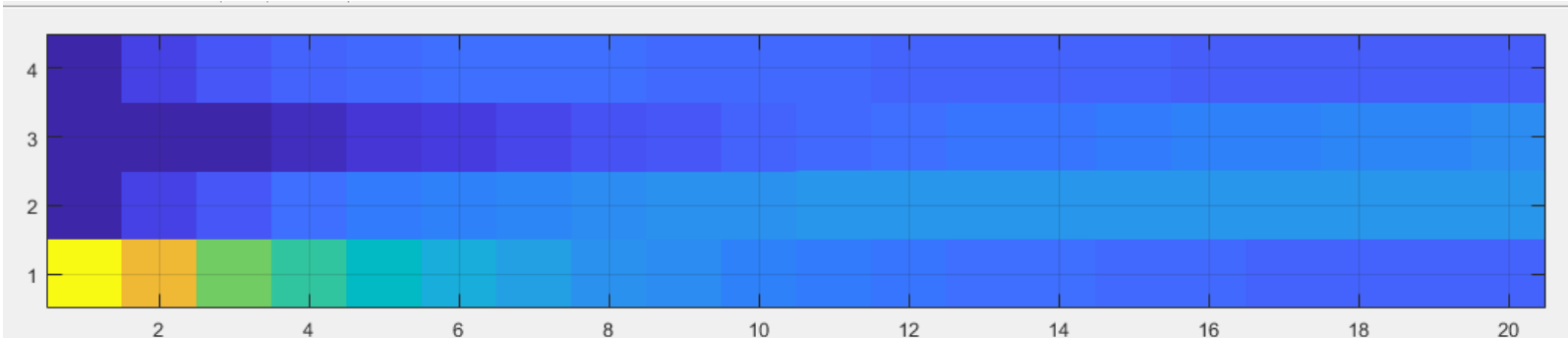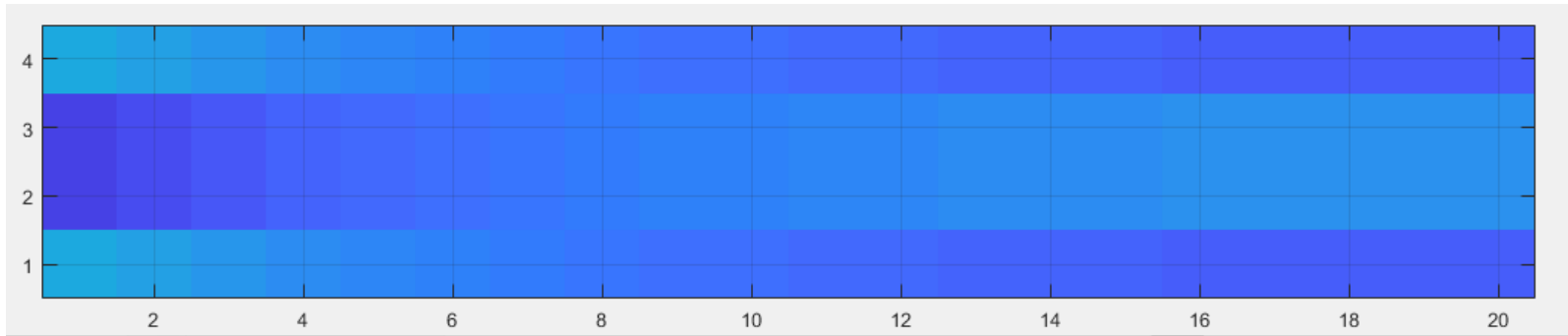


No return

# Limiting distribution

# When the Markov chain is irreducible / aperiodic /recurrent

- The limiting distribution is equal to the stationary distribution
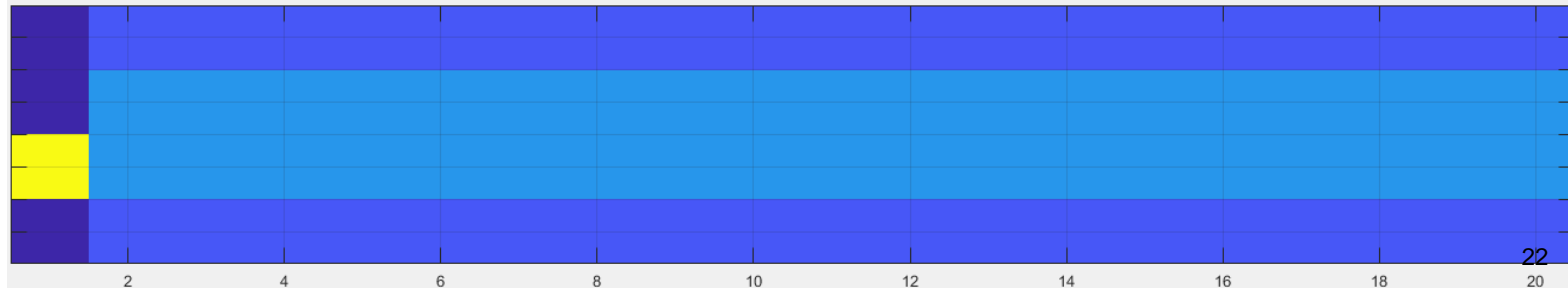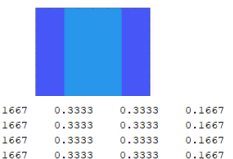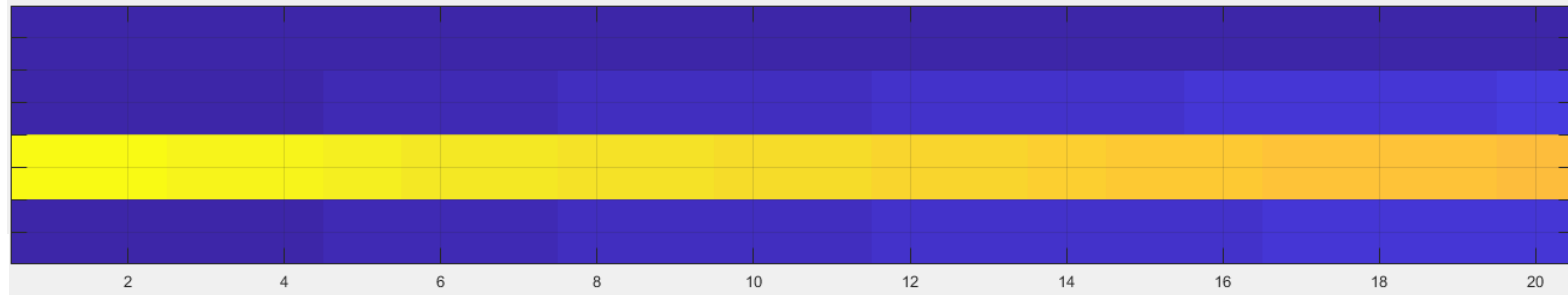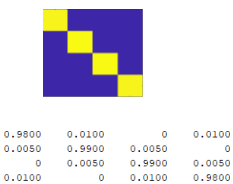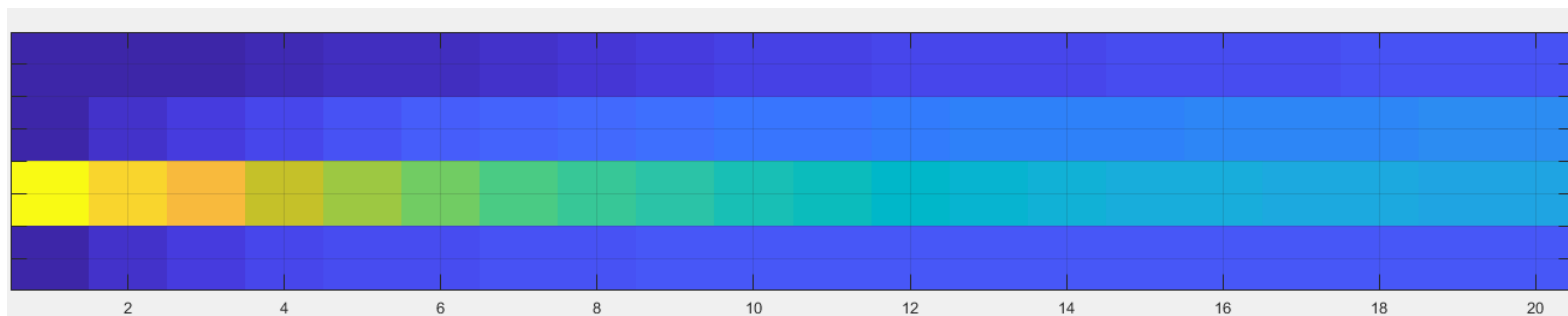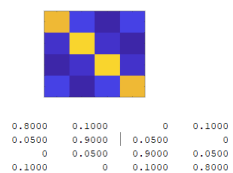
$$p_s = p_{\text{Lim}}$$
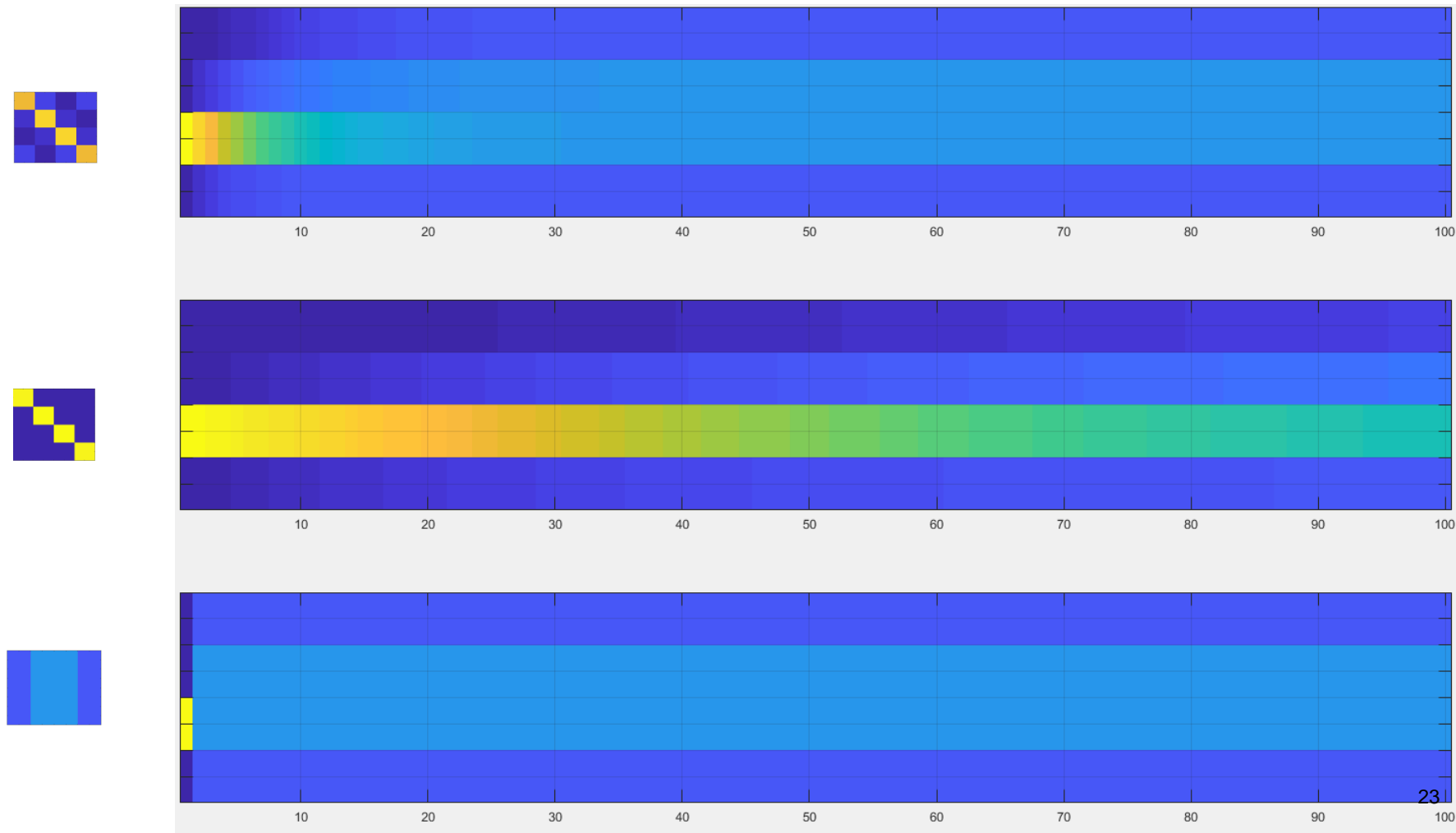
- Stationary distribution is fix point of iteration

$$p_s P = p_s$$

- Limiting distribution (is independent of $p_0$ )

$$\lim_{n \to \infty} p_0 P^n = p_{\text{Lim}}$$

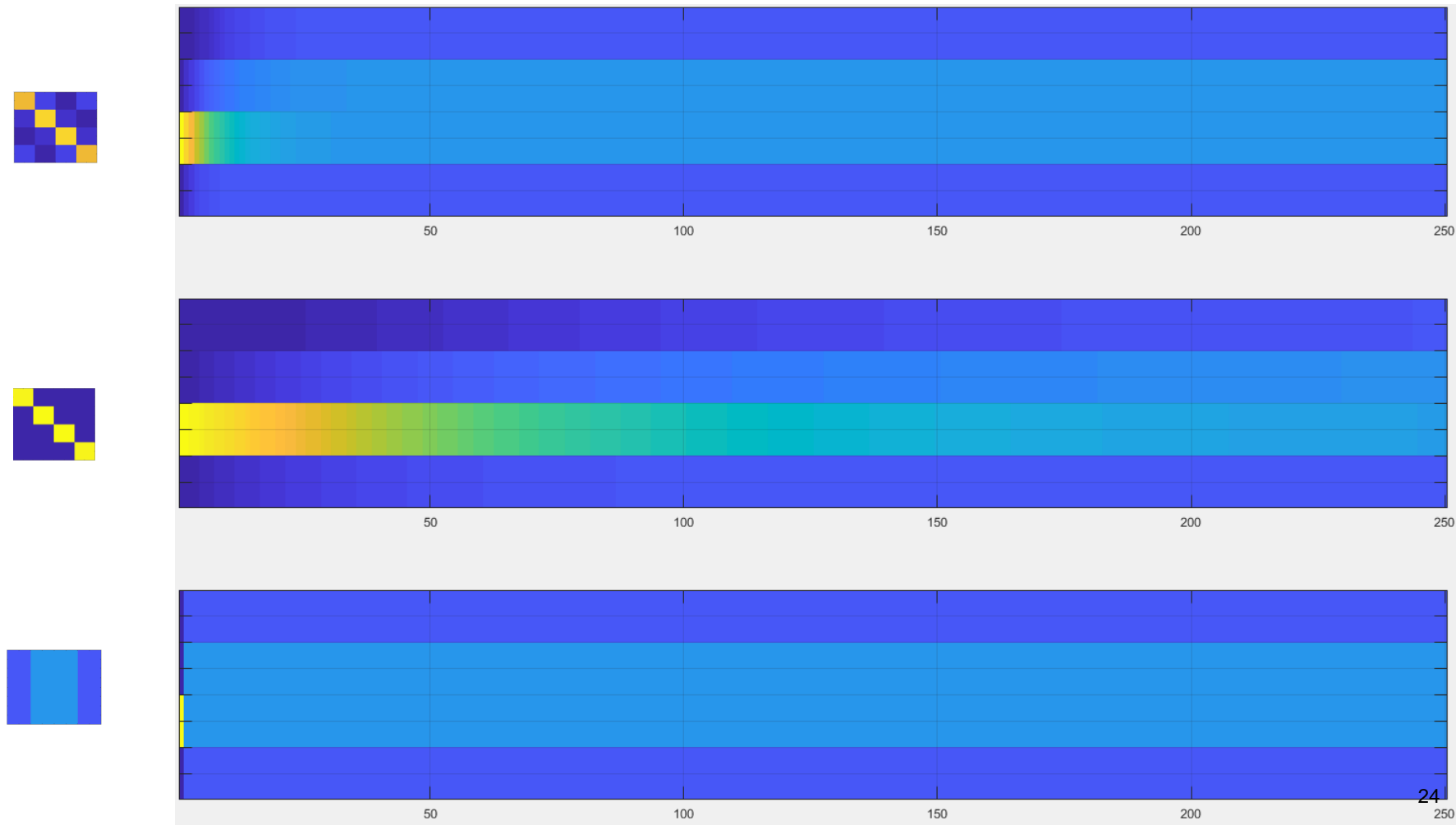# Time to reach limiting distribution n=20

# Time to reach limiting distribution n=100

# Time to reach limiting distribution n=250

# Markov chain theory – discrete case

- Assume $\{X^{(t)}\}$ is a Markov chain where $X^{(t)}$ is a discrete random variable

$$\Pr(X^{(t)} = y | X^{(t-1)} = x) = P(y|x)$$

  giving the transition probabilities
- Assume the chain is
  - irreducible: It is possible to move from any **x** to any **y** in a finite number of steps
  - reccurent: The chain will visit any state infinitely often.
  - aperiodic: Does not go in cycles
- Then there exists a unique distribution $f(x)$ such that

$$\lim_{t \to \infty} \Pr(X^{(t)} = y | X^{(0)} = x) = f(y)$$

$$\hat{\mu}_{MCMC} \to \mu = E^f[X]$$

- How to find $f(\cdot)$ (the stationary distribution): Solve

$$f(y) = \sum_x f(x) P(y|x)$$

- Our situation: We have $f(y)$, want to find $P(y|x)$
  - Note: Many possible $P(y|x)$!

25

# Markov chain theory - general setting

- Assume $\{\mathbf{X}^{(t)}\}$ is a Markov chain where $\mathbf{X}^{(t)} \in S$

$$\Pr(\mathbf{X}^{(t)} \in A | \mathbf{X}^{(t-1)} = \mathbf{x}) = P(\mathbf{x}, A) = \int_{\mathbf{y} \in A} P(\mathbf{y}|\mathbf{x}) d\mathbf{y}$$

giving the transition densities
- Assume the chain is
  - irreducible: It is possible to move from any $\mathbf{x}$ to any $\mathbf{y}$ in a finite number of steps
  - reccurent: The chain will visit any $A \subset S$ infinitely often.
  - aperiodic: Do not go in cycles
- Then there exists a distribution $f(\mathbf{x})$ such that

$$\lim_{t \to \infty} \Pr(\mathbf{X}^{(t)} \in A | \mathbf{X}^{(0)} = \mathbf{x}) = \int_A f(\mathbf{y}) d\mathbf{y}$$
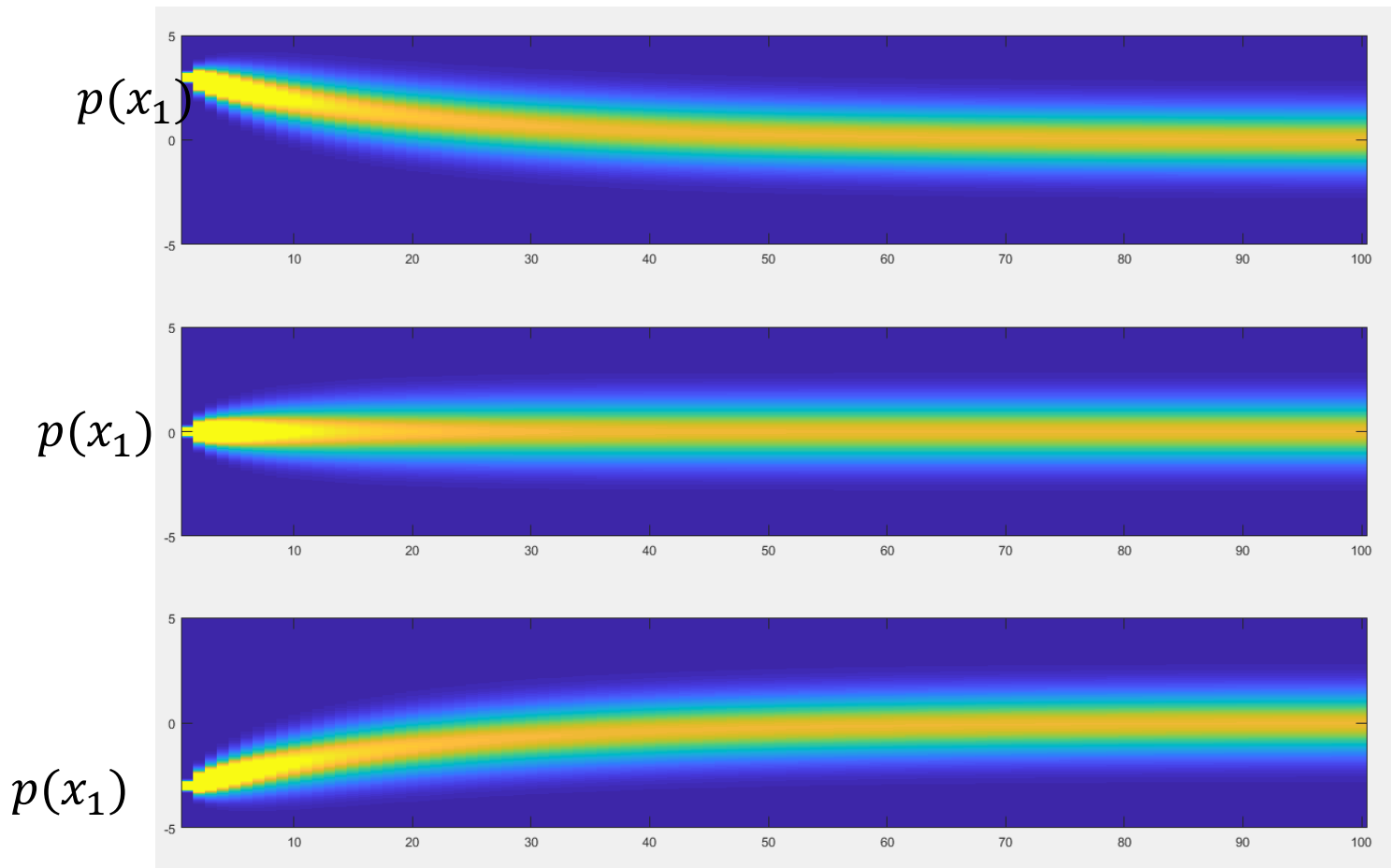
$$\hat{\mu}_{MCMC} \to \mu$$

- How to find $f(\cdot)$ (the stationary distribution): Solve

$$f(\mathbf{y}) = \int_{\mathbf{x}} f(\mathbf{x}) P(\mathbf{y}|\mathbf{x}) d\mathbf{x}$$

- Our situation: We have $f(\cdot)$, want to find $P(\mathbf{y}|\mathbf{x})$

# Example of a continuous transition density, AR1 model

$$p(x_t | x_{t-1}) = \phi(ax_{t-1}, \sigma^2(1 - a^2))$$

$p(x_1)$

$p(x_1)$

$p(x_1)$



27

# Questions

# We want to construct P(x|y) to match our needs

- Need to have good properties
  - Stationary
  - Irreducible
  - Aperiodic
  - Recurrent

- Also need to get our target as a stationary distribution

$$f(\mathbf{y}) = \int_{\mathbf{x}} f(\mathbf{x}) P(\mathbf{y}|\mathbf{x}) d\mathbf{x}$$

  - Simplify the hunt by introducing symmertry
  - **detailed balance**

# Detailed balance

- The task: Find a transition probability/density $P(\mathbf{y}|\mathbf{x})$ satisfying

$$f(\mathbf{y}) = \int_{\mathbf{x}} f(\mathbf{x})P(\mathbf{y}|\mathbf{x})d\mathbf{x}$$

Can in general be a difficult criterion to check
- Sufficient criterion:

$$f(\mathbf{x})P(\mathbf{y}|\mathbf{x}) = f(\mathbf{y})P(\mathbf{x}|\mathbf{y}) \quad \text{Detailed balance}$$

We then have

$$\int_{\mathbf{x}} f(\mathbf{x})P(\mathbf{y}|\mathbf{x})d\mathbf{x} = \int_{\mathbf{x}} f(\mathbf{y})P(\mathbf{x}|\mathbf{y})d\mathbf{x}$$

$$= f(\mathbf{y})\int_{\mathbf{x}} P(\mathbf{x}|\mathbf{y})d\mathbf{x} = f(\mathbf{y})$$

since $P(\mathbf{x}|\mathbf{y})$ is, for any given $\mathbf{y}$, a density wrt $\mathbf{x}$.
- Note: For $\mathbf{y} = \mathbf{x}$, detailed balance always fulfilled, only necessary to check for $\mathbf{y} \neq \mathbf{x}$.

# Metropolis-Hastings algorithm

- $P(\mathbf{y}|\mathbf{x})$ defined through an algorithm:
  1. Sample a candidate value $\mathbf{X}^*$ from a proposal distribution $g(\cdot|\mathbf{x})$.
  2. Compute the Metropolis-Hastings ratio

  $$R(\mathbf{x}, \mathbf{X}^*) = \frac{f(\mathbf{X}^*)g(\mathbf{x}|\mathbf{X}^*)}{f(\mathbf{x})g(\mathbf{X}^*|\mathbf{x})}$$

  3. Put

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

- For $\mathbf{y} \neq \mathbf{x}$:

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

- Note: $P(\mathbf{x}|\mathbf{x})$ somewhat difficult to evaluate in this case.

Either we keep **x** with a certain probability
Or we change to **X*** which have a certain density

# Metropolis-Hastings algorithm
# Detailed balance

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

$$= \min\{f(\mathbf{x})g(\mathbf{y}|\mathbf{x}), f(\mathbf{y})g(\mathbf{x}|\mathbf{y})\}$$

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

# The probability of a value being repeated is positive

Pf:

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

$$\int_{y\neq x} P(y|x)dy = \int_{y\neq x} \underbrace{g(y|x)}_{} \underbrace{\min\left\{1, \frac{f(y)g(x|y)}{f(x)g(y|x)}\right\}}_{} dy \leq 1$$

Density:
integrates to 1

Positive number:
$\leq 1$

# What about unknown scaling and MH

- Assume now $f(\mathbf{x}) = c \cdot q(\mathbf{x})$ with $c$ unknown.

$$R(\mathbf{x},\mathbf{y}) = \frac{f(\mathbf{y})g(\mathbf{x}|\mathbf{y})}{f(\mathbf{x})g(\mathbf{y}|\mathbf{x})} = \frac{c \cdot q(\mathbf{y})g(\mathbf{x}|\mathbf{y})}{c \cdot q(\mathbf{x})g(\mathbf{y}|\mathbf{x})} = \frac{q(\mathbf{y})g(\mathbf{x}|\mathbf{y})}{q(\mathbf{x})g(\mathbf{y}|\mathbf{x})}$$

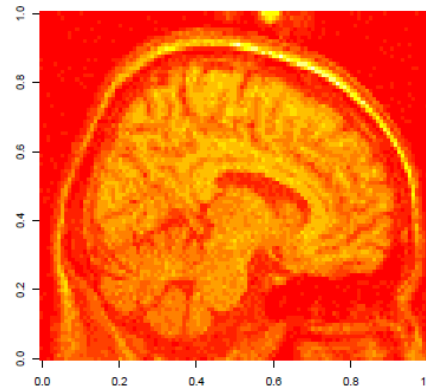

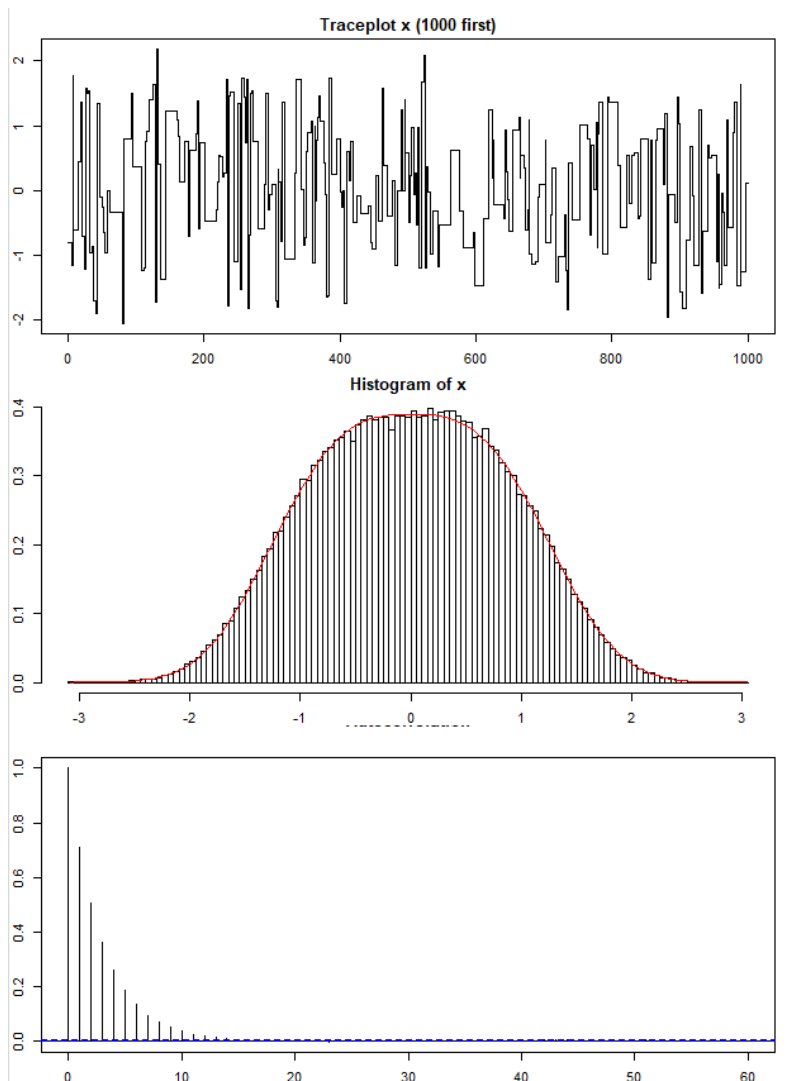

- Do not depend on $c$!

Important for Bayesian analysis Posterior ∝ Likelihood × Prior

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)} \propto p(y|x)p(x)$$

**Important for Gibbs type distributions**

$$\Pr(\mathbf{C}) = \Pr(C_{11}, \ldots., C_{n_1 n_2})$$
$$= \frac{1}{Z} e^{-\beta \sum_{||(i,j)-(i'j')||=1} I(C_{ij} \neq C_{i'j'})}$$

$$\Pr(\mathbf{C}|\mathbf{y}) = \frac{\Pr(\mathbf{C}) \prod_{ij} f(y_{ij}|C_{ij})}{\sum_{\mathbf{C}'} \Pr(\mathbf{C}') \prod_{ij} f(y_{ij}|C'_{ij})}$$

# Questions

# **Metropolis Hastings is a general form:**

- Specific chains:
  - Random walk chains
  - Independent chains
  - Gibbs sampler
- Tricks to customize sampling
  - Reparametrize
  - Block update
  - Hybrid
  - Griddy Gibbs

# Random walk chains

- Popular choice of proposal distribution:

$$\mathbf{X}^* = \mathbf{x} + \boldsymbol{\varepsilon}$$

- $g(\mathbf{x}^*|\mathbf{x}) = h(\mathbf{x}^* - \mathbf{x})$
- Popular choices: Uniform, Gaussian, $t$-distribution
- Note: If $h(\cdot)$ is symmetric, $g(\mathbf{x}^*|\mathbf{x}) = g(\mathbf{x}|\mathbf{x}^*)$ and

$$R(\mathbf{x}, \mathbf{x}^*) = \frac{f(\mathbf{x}^*)g(\mathbf{x}|\mathbf{x}^*)}{f(\mathbf{x})g(\mathbf{x}^*|\mathbf{x})} = \frac{f(\mathbf{x}^*)}{f(\mathbf{x})}$$

# Example

- Assume $f(x) \propto \exp(-|x|^3/3)$

- Proposal distribution $N(x, 4^2)$

- Example_MH_cubic.R

```
#Initial value
x = rnorm(1)
acc = 0
for(i in 2:N)
{
  y = rnorm(1,x[i-1],4) # proposal
  R = f(y)/f(x[i-1])    # acceptance ratio
                        # note that the acceptance rate is min(1,R),
  if(runif(1)<R)        # The syntax her will give that since we allways accept if R>1
  {
    x[i] = y
    acc = acc+1
  }
  else
   x[i] = x[i-1]
}
```

# Results random walk



Acceptance rate
= 0.2755276

Lag one scatterplot

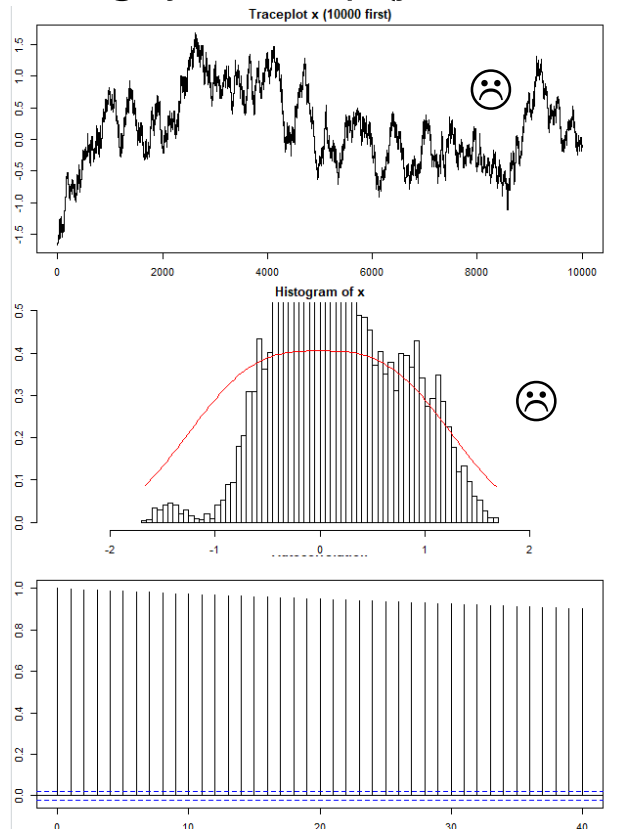# The repeats of a value is needed to get the correct distribution

Compare histograms to true distribution

This is kind of similar to what we have for sampling importance resampling (SIR) If a value is repeated it gets «more weight»
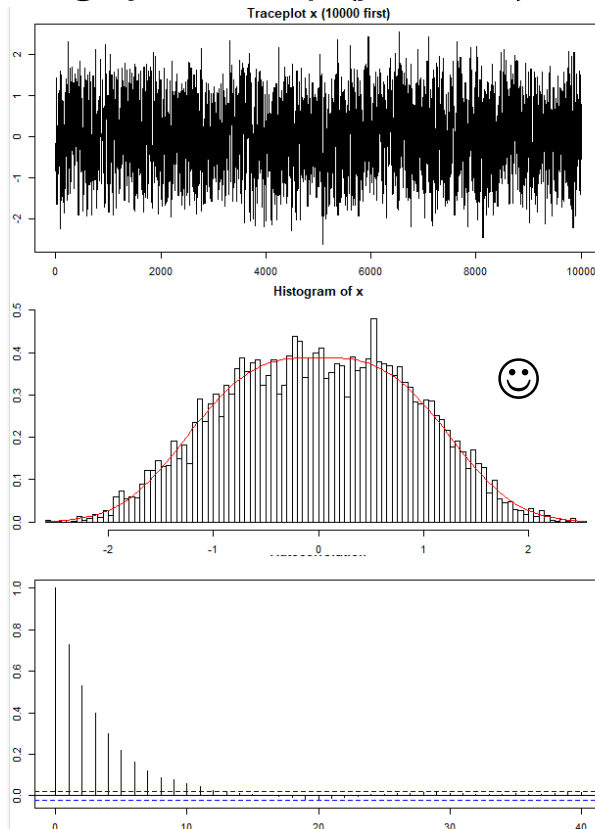


40

# The effect variance in proposal distribution

$$g(y|x) = \phi(y; x, 0.04^2) \qquad g(y|x) = \phi(y; x, 1^2) \qquad g(y|x) = \phi(y; x, 100^2)$$



Acc. rate = 0.994

Too small steps,
high acceptance
high correlation ☹

Acc. Rate = 0.700

Just about right,
good acceptance
low correlation ☺

Acc. Rate = 0.012

Too large changes proposed,
low acceptance
high correlation ☹

41

# Questions?

# Independent chains

- Assume $g(\mathbf{x}^*|\mathbf{x}) = g(\mathbf{x}^*)$. Then

$$R(\mathbf{x}, \mathbf{x}^*) = \frac{f(\mathbf{x}^*)g(\mathbf{x})}{f(\mathbf{x})g(\mathbf{x}^*)} = \frac{\frac{f(\mathbf{x}^*)}{g(\mathbf{x}^*)}}{\frac{f(\mathbf{x})}{g(\mathbf{x})}},$$

fraction of <span style="color:red">importance weights</span>!

- Behave very much like importance sampling and SIR
- Difficult to specify $g(\mathbf{x})$ for high-dimensional problems
- Theoretical properties easier to evaluate than for random walk versions.

Challenges similar to what seen in:
- rejection sampling
- importance sampling
- sampling importance resampling

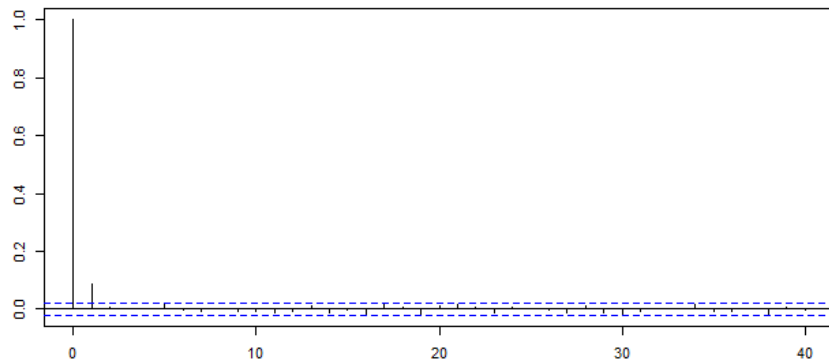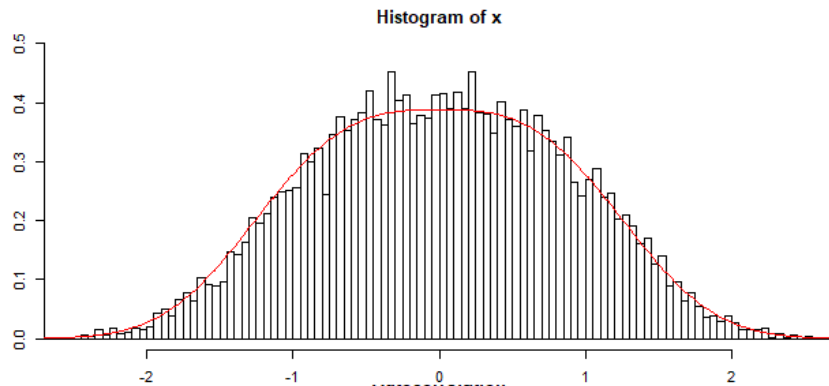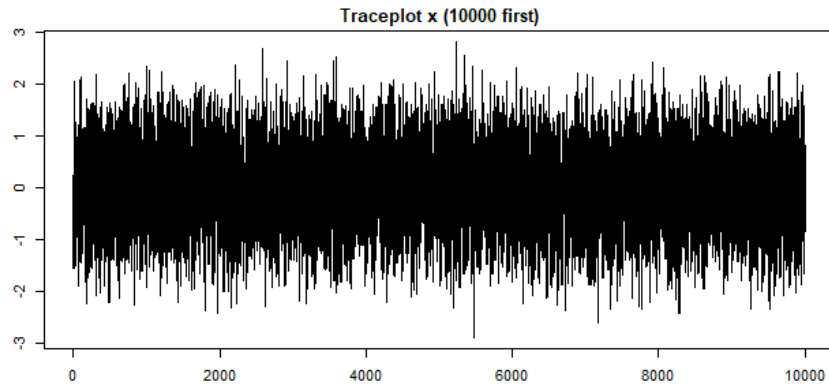# Example

- Assume $f(x) \propto \exp(-|x|^3/3)$

$$g(y|x) = \phi(y; 0, 1^2)$$

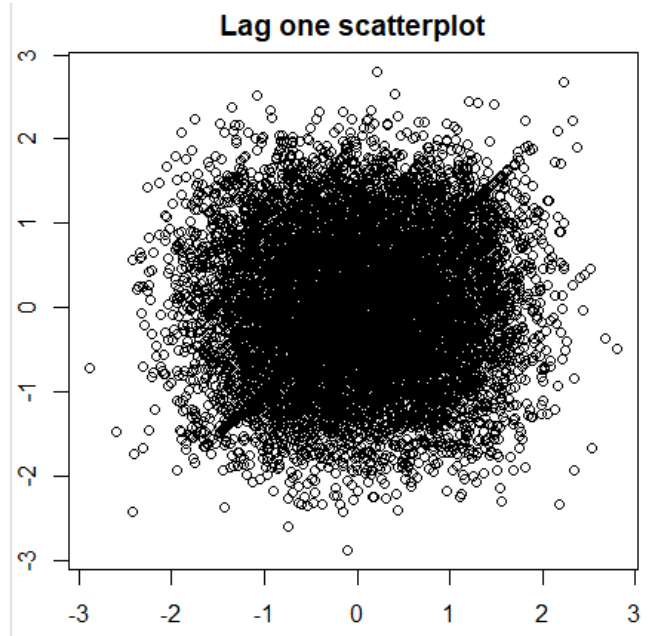`Example_MH_cubic_independence.R`

```
N = 10000      # Number of iterations
x = rep(NA,N)
varProp=1^2 # variance of proposal

#Initial value
x =  rnorm(1,0,varProp)
acc = 0
for(i in 2:N)
{
  y = rnorm(1,0,varProp) # proposal
  R = f(y)*dnorm(x[i-1],0,varProp)/(f(x[i-1])*dnorm(y,0,varProp))          # acceptance ratio
                         # note that the acceptance rate is min(1,R),
  if(runif(1)<R)         # The syntax her will give that since we allways accept if R>1
  {
    x[i] = y
    acc = acc+1
  }
  else
   x[i] = x[i-1]
}
```
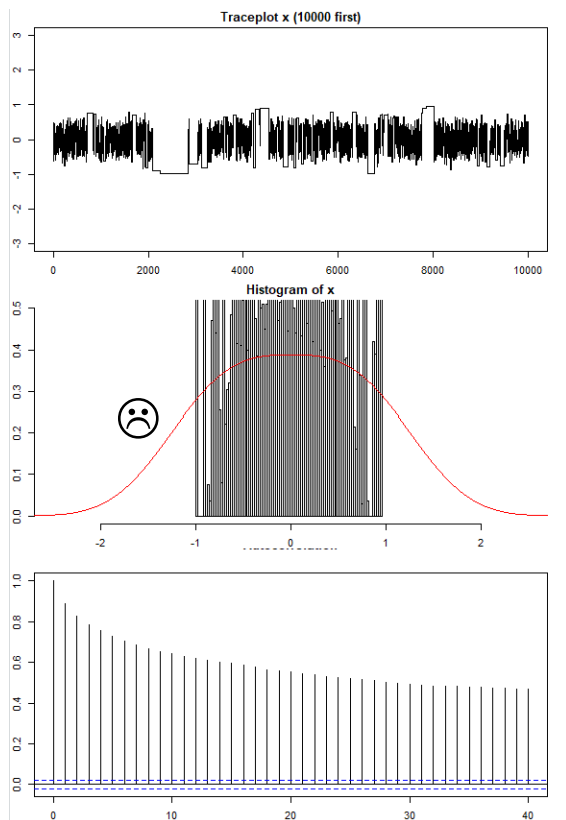
# Results independent



Acceptance rate= 0.9149915

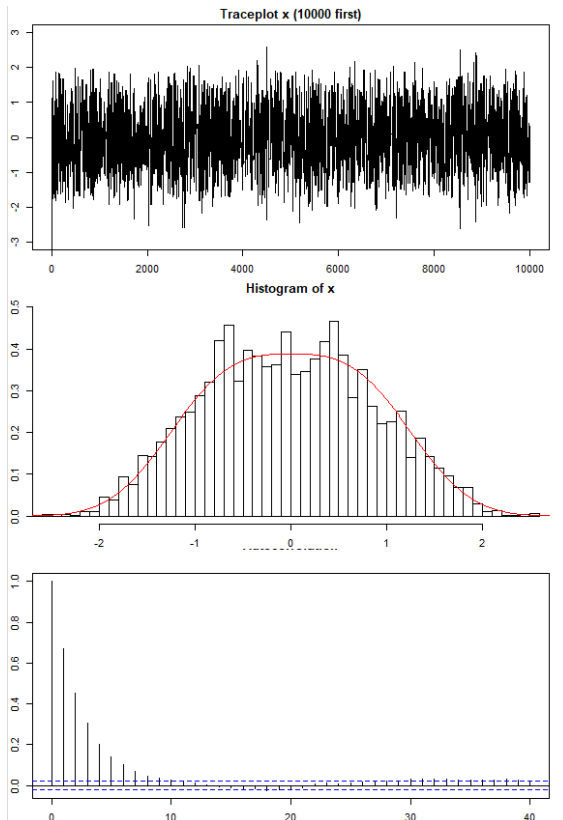# The effect variance in proposal distribution

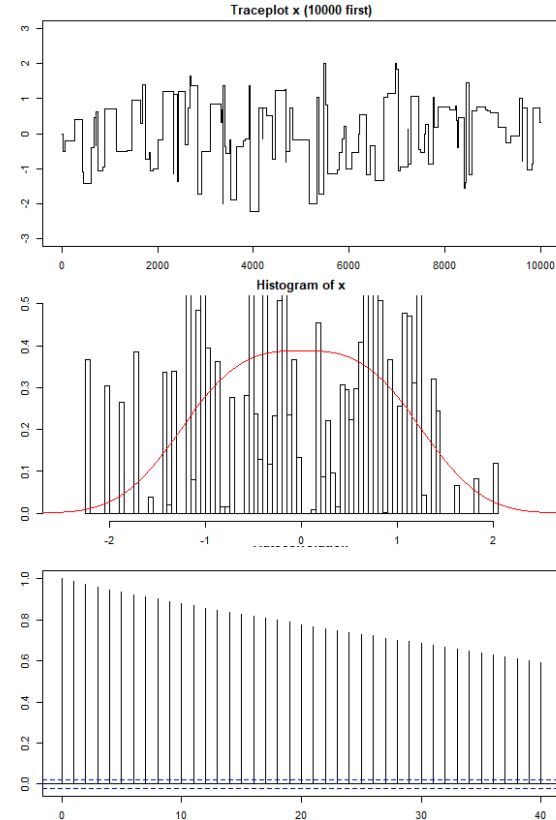$$g(y|x) = \phi(y; 0, 0.25^2) \qquad g(y|x) = \phi(y; 0, 4^2) \qquad g(y|x) = \phi(y; 0, 100^2)$$



Acc. rate =0.419

Too narrow proposal,
good acceptance
high correlation ☹

Acc. Rate = 0.288

Just about right,
reasonable acceptance
low correlation ☺

Acc. Rate = 0.012

Too large changes proposed,
low acceptance
high correlation ☹

46

# Questions

# M-H and multivariate settings

- $\mathbf{X} = (X_1, ..., X_p)$
- Typical in this case: Only change <span style="color:red">one</span> or a few components at a time.
  1. Choose index $j$ (randomly)
  2. Sample $X_j^* \sim g_j(\cdot|\mathbf{x})$, put $X_k^* = X_k$ for $k \neq j$
  3. Compute

$$R(\mathbf{x}, \mathbf{X}^*) = \frac{f(\mathbf{X}^*)g(\mathbf{x}|\mathbf{X}^*)}{f(\mathbf{x})g(\mathbf{X}^*|\mathbf{x})}$$

  4. Put

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

- Can show that this version also satisfies detailed balance
- Can even go through indexes systematic
  - Should then consider the whole loop through all components as one iteration

# Example multivariate with single coordinate update

- Assume $f(\mathbf{x}) \propto \exp(-||\mathbf{x}||^3/3) = \exp(-[||\mathbf{x}||^2]^{3/2}/3)$
- Proposal distribution
  1. $j \sim \text{Uniform}[1, 2, ..., p]$
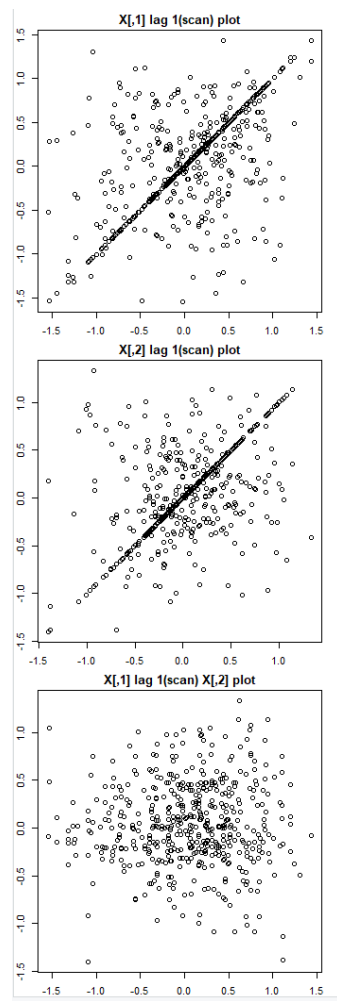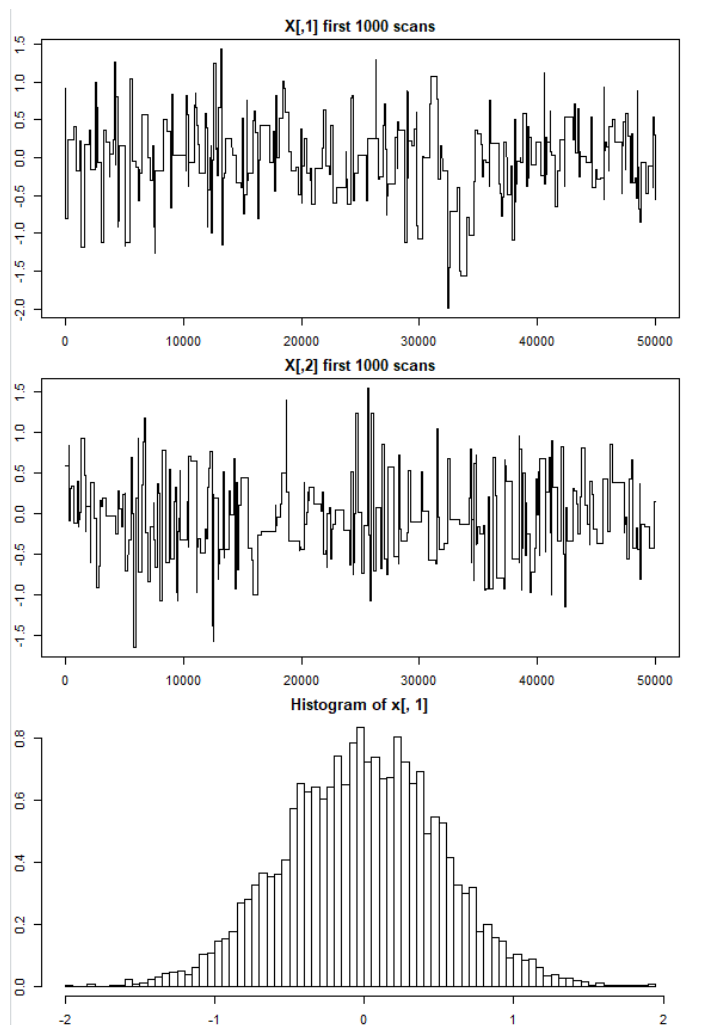  2. $x_j^* \sim N(x_j, 1)$
- Example_MH_cubic_multivariate.R

```
#Proposal distribution: Gaussian distribution centered at previous value
p = 50
N = 10000     # Number of iterations
x = matrix(nrow=N,ncol=p)

#Initial value
x[1,] = rnorm(p)
acc = 0
for(i in 2:N)
{
  j = sample(1:p,1)
  y = x[i-1,]
  y[j] = rnorm(1,x[i-1,j],2)
  R = f(y)*dnorm(x[i-1,j],y[j],1)/(f(x[i-1,])*dnorm(y[j],x[i-1,j],1))
  if(runif(1)<R)
  {
    x[i,] = y
    acc = acc+1
  }
  else
   x[i,] = x[i-1,]
}
```

See also fixed scan in: `Example_MH_cubic_multivariate_2.R`

# Results independent



Acceptance rate= 0.3053943

50

# Questions