

Statistical Techniques for Data Science & Robotics

Week 11

Objectives for today

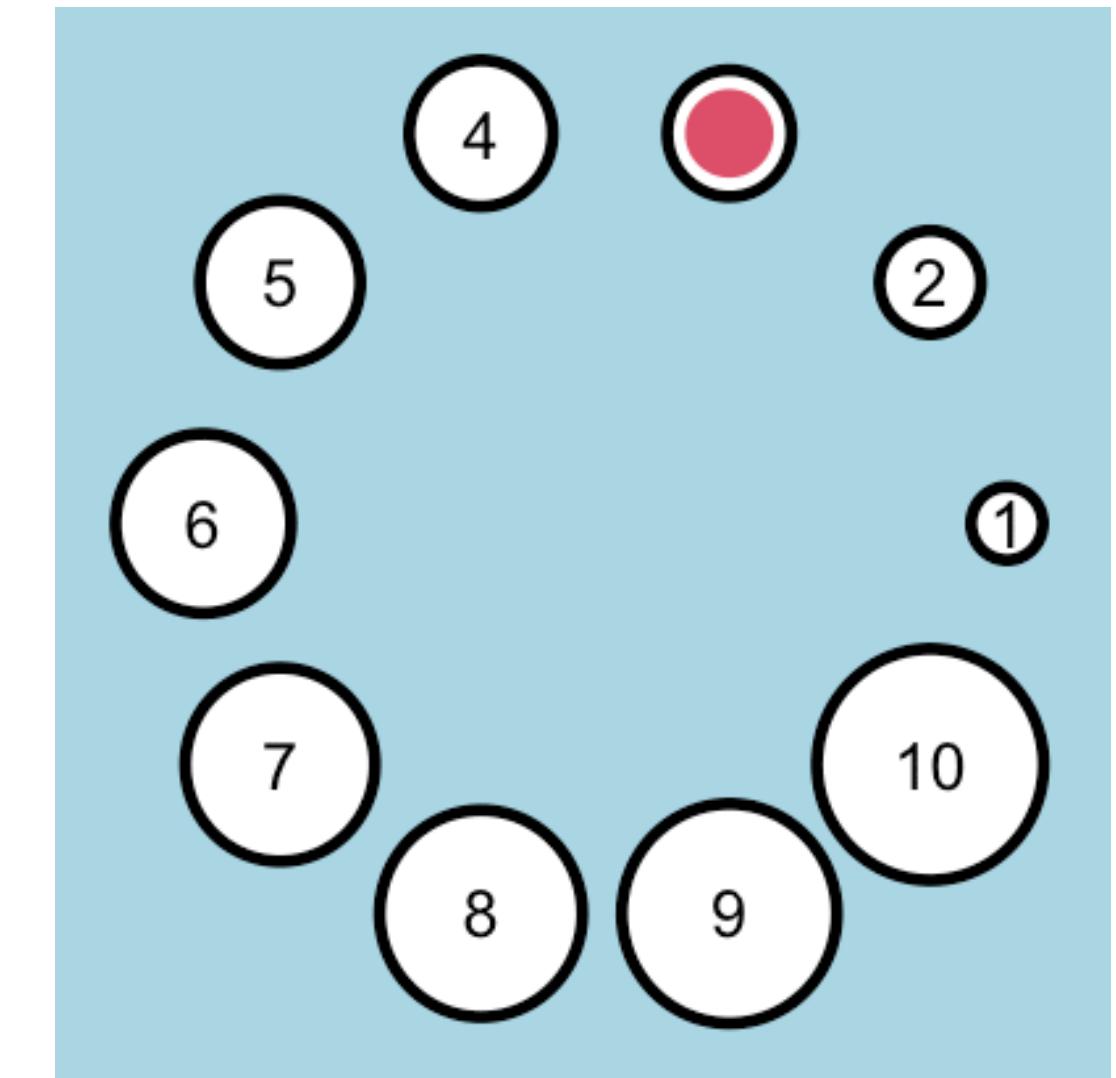
- Markov Chains
- Metropolis-Hastings algorithm
 - Random walk Metropolis
 - Metropolis-Hastings



King Markov

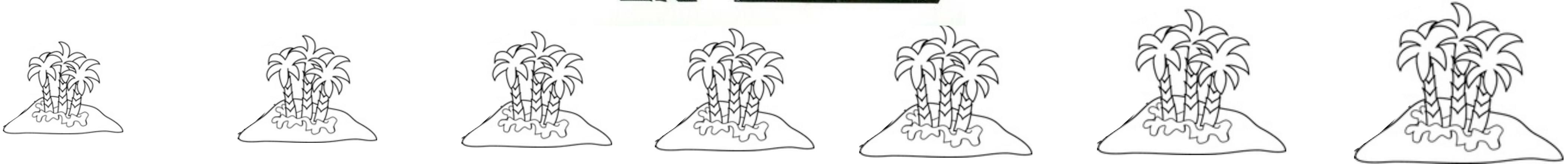
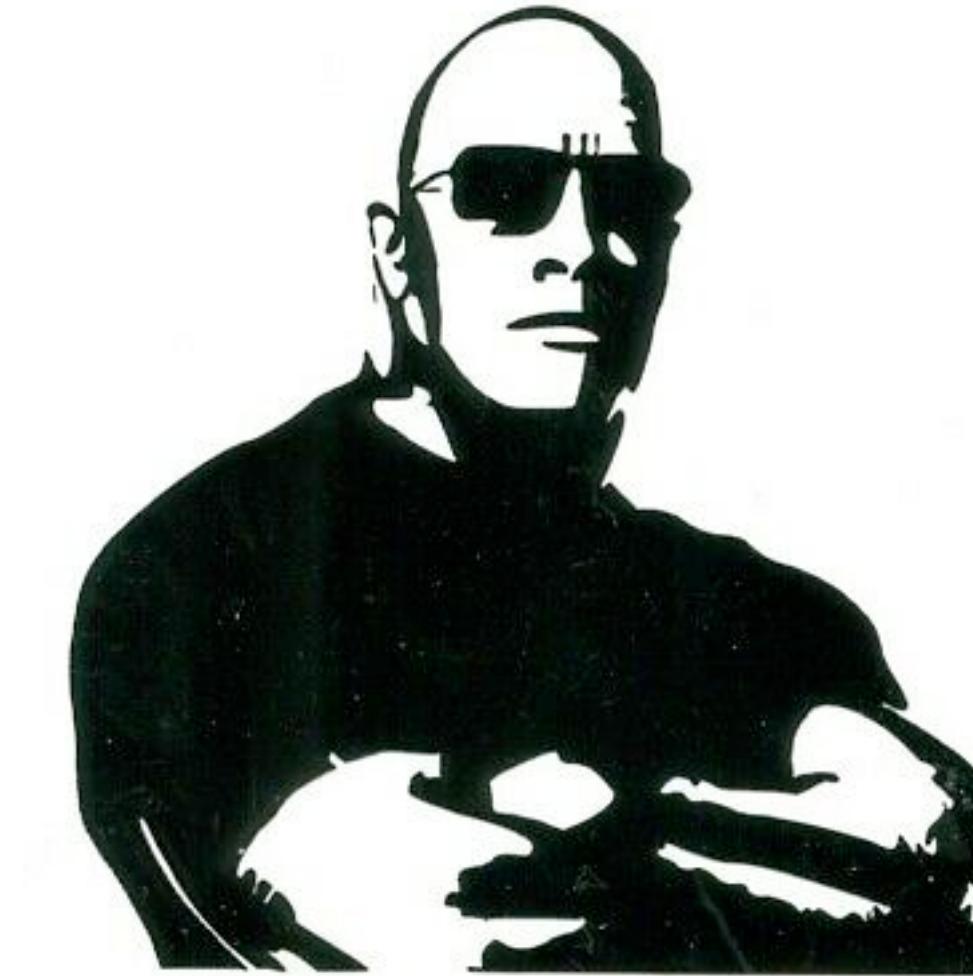
- King Markov was a benevolent autocrat of an island kingdom, a circular archipelago, with 10 islands.

Source: <https://speakerdeck.com/rmcelreath/statistical-rethinking-2023-lecture-08>

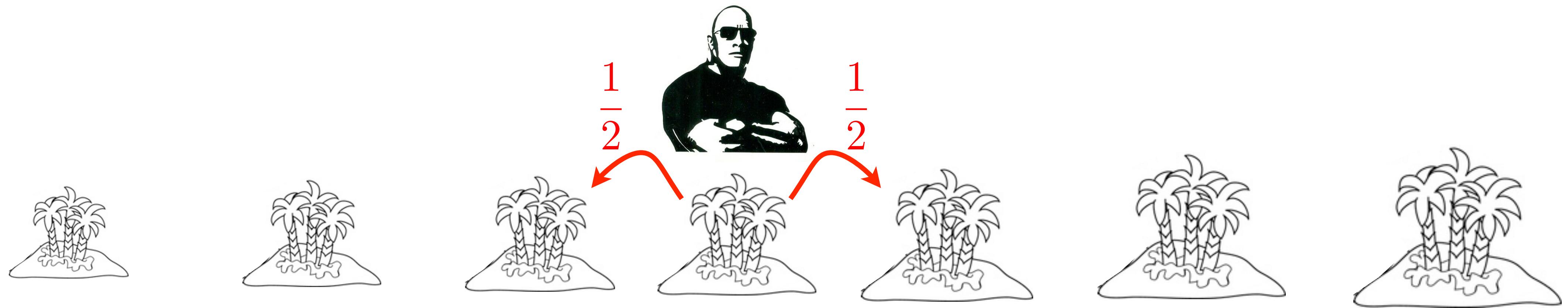


The Metropolis Archipelago

Contract: King Markov must visit each island in proportion to its population size

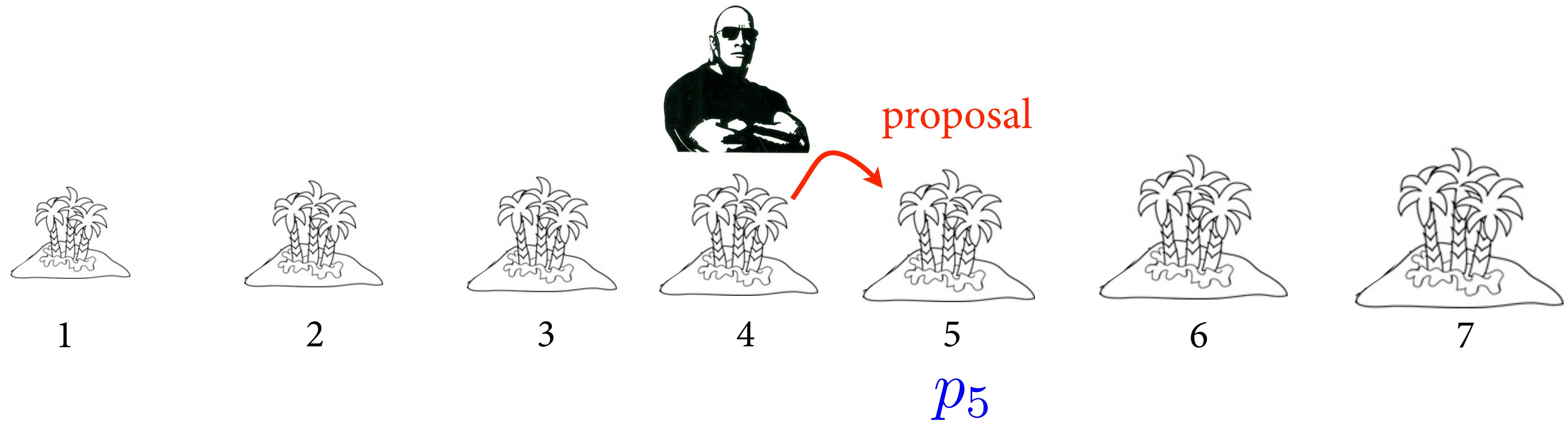


Here's how he does it...



(1) Flip a coin to choose island on left or right.
Call it the “proposal” island.

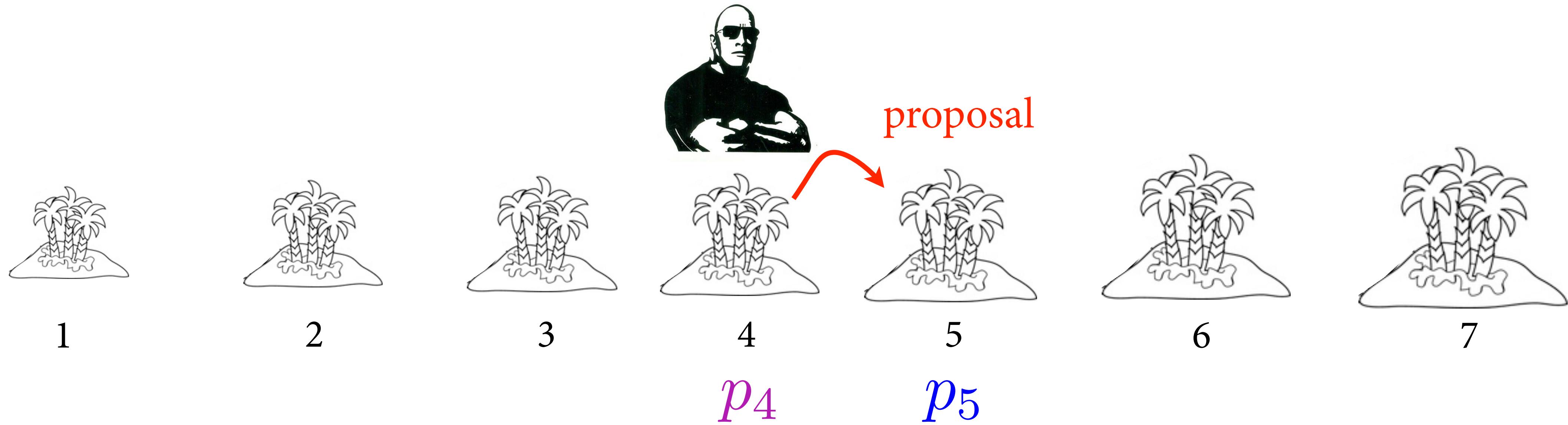
(1) Flip a coin to choose island on left or right.
Call it the “proposal” island.



(2) Find population of proposal island.

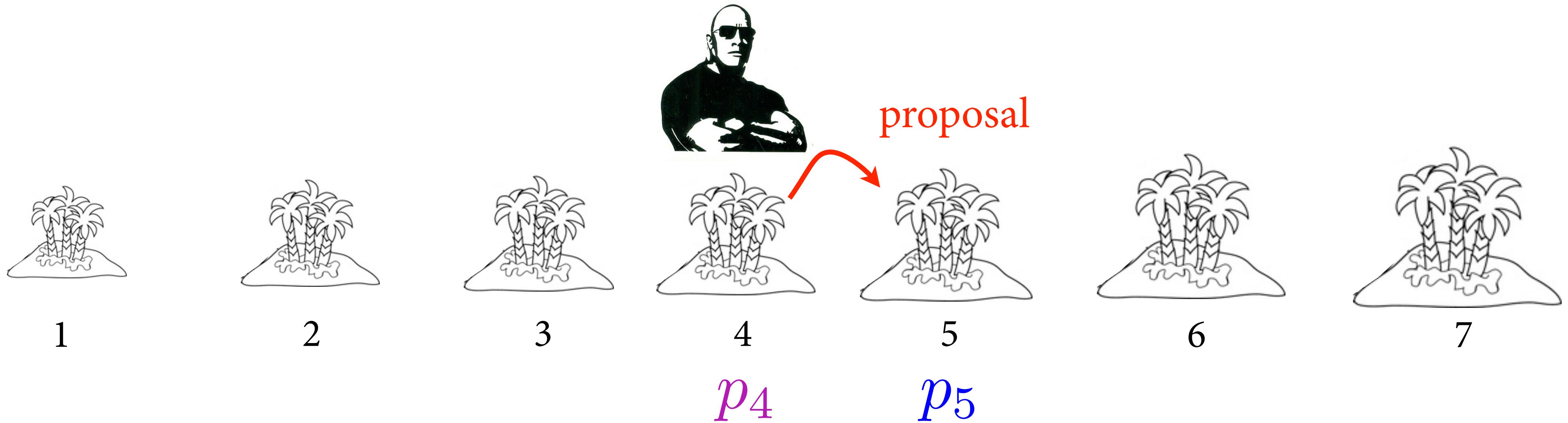
(1) Flip a coin to choose island on left or right.
Call it the “proposal” island.

(2) Find population of proposal island.



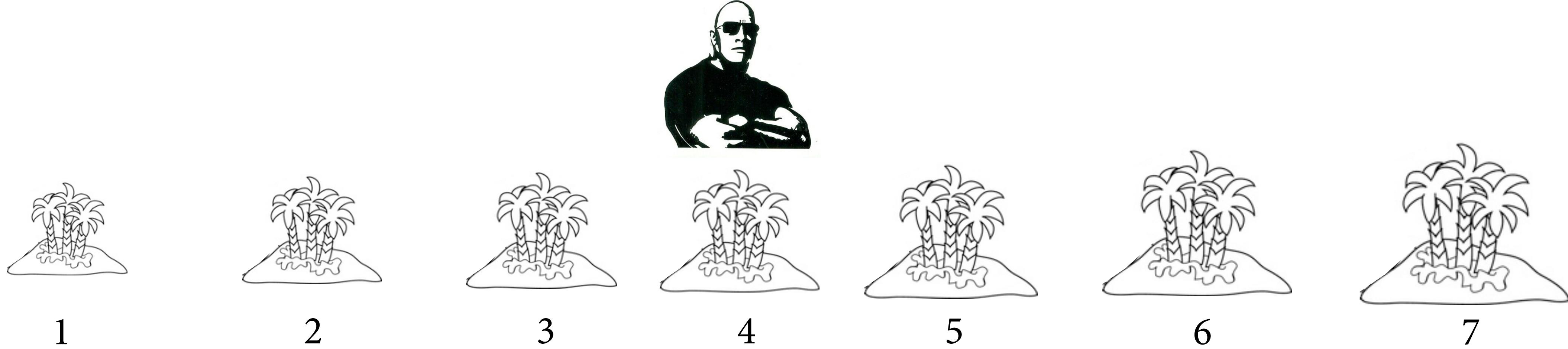
(3) Find population of current island.

- (1) Flip a coin to choose island on left or right.
Call it the “proposal” island.
- (2) Find population of proposal island.
- (3) Find population of current island.



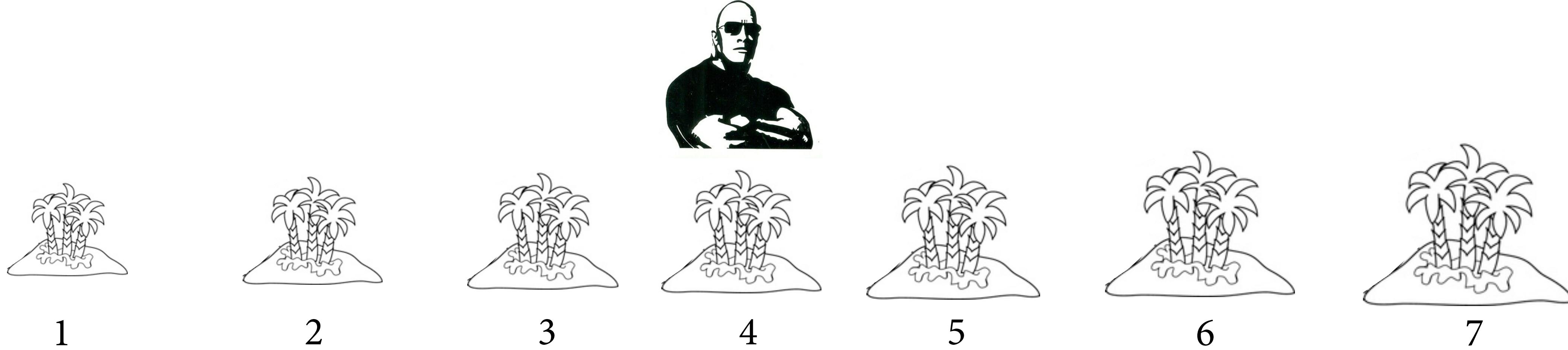
(4) Move to proposal, with probability = $\frac{p_5}{p_4}$

- (1) Flip a coin to choose island on left or right.
Call it the “proposal” island.
- (2) Find population of proposal island.
- (3) Find population of current island.
- (4) Move to proposal, with probability = $\frac{p_5}{p_4}$

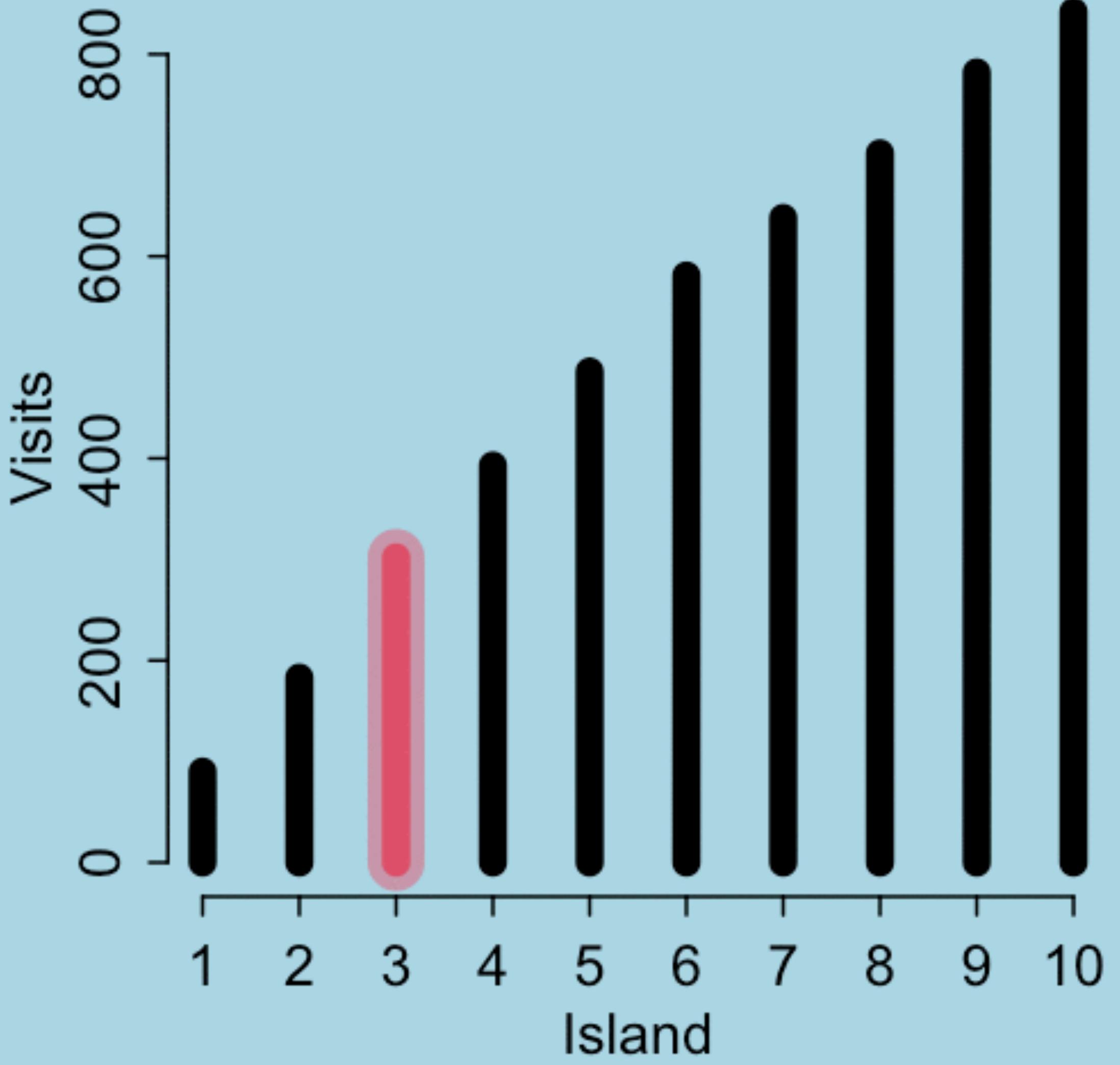
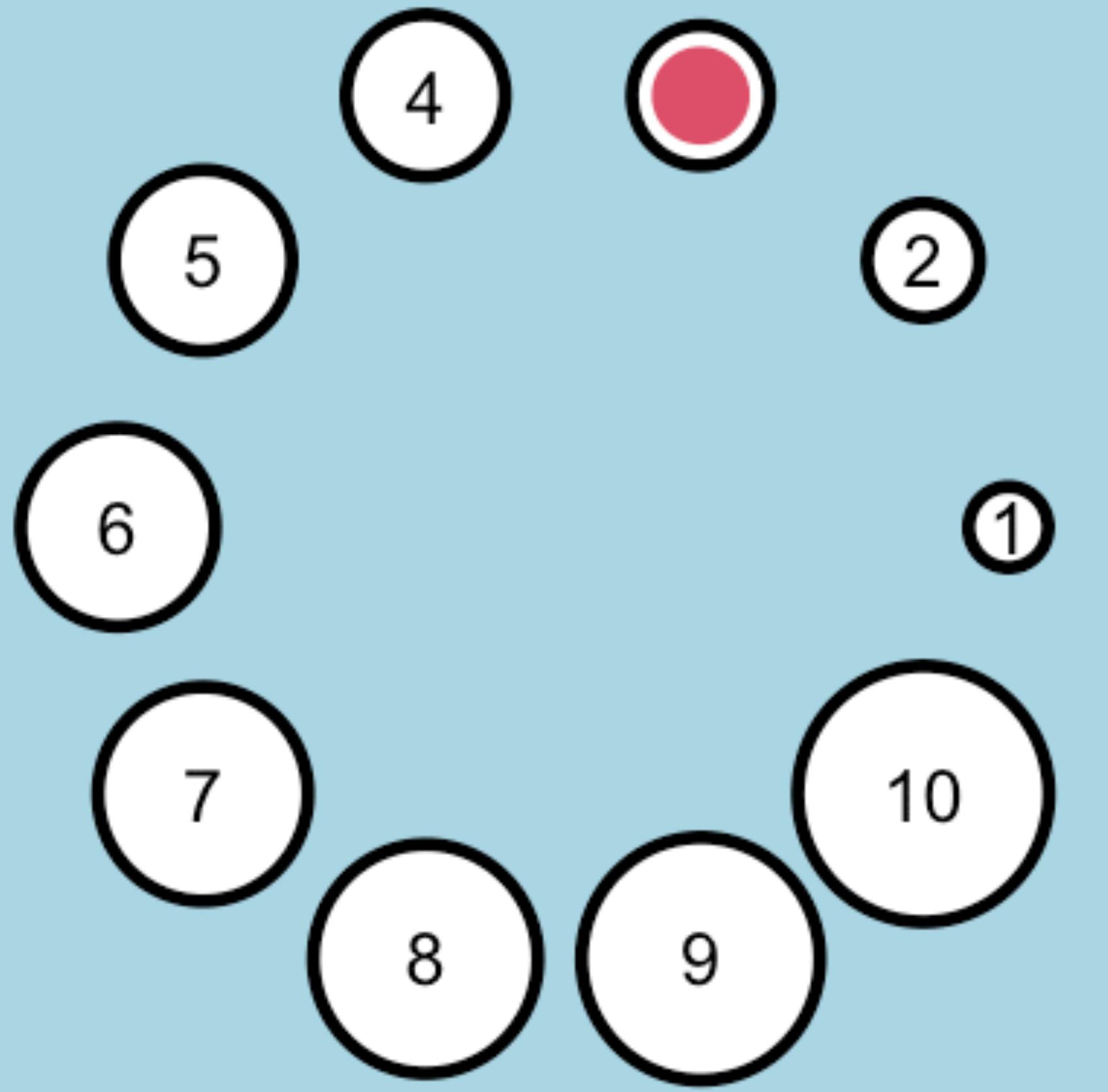


(5) Repeat from (1)

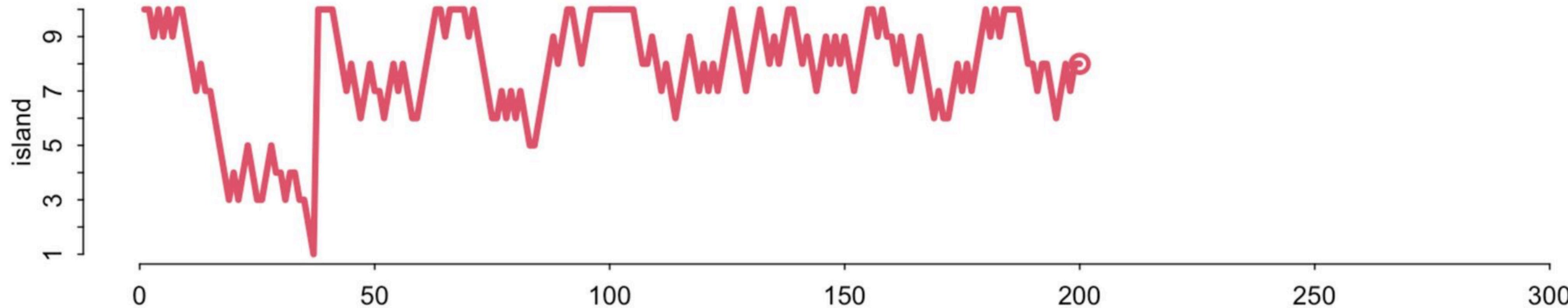
- (1) Flip a coin to choose island on left or right.
Call it the “proposal” island.
- (2) Find population of proposal island.
- (3) Find population of current island.
- (4) Move to proposal, with probability = $\frac{p_5}{p_4}$
- (5) Repeat from (1)



This procedure ensures visiting each island in proportion to its population, *in the long run*.



“Markov chain Monte Carlo”



Metropolis algorithm: Simple version of *Markov chain Monte Carlo* (MCMC)

Easy to write, very general, often inefficient

A tale of Good King Markov



Markov Chains

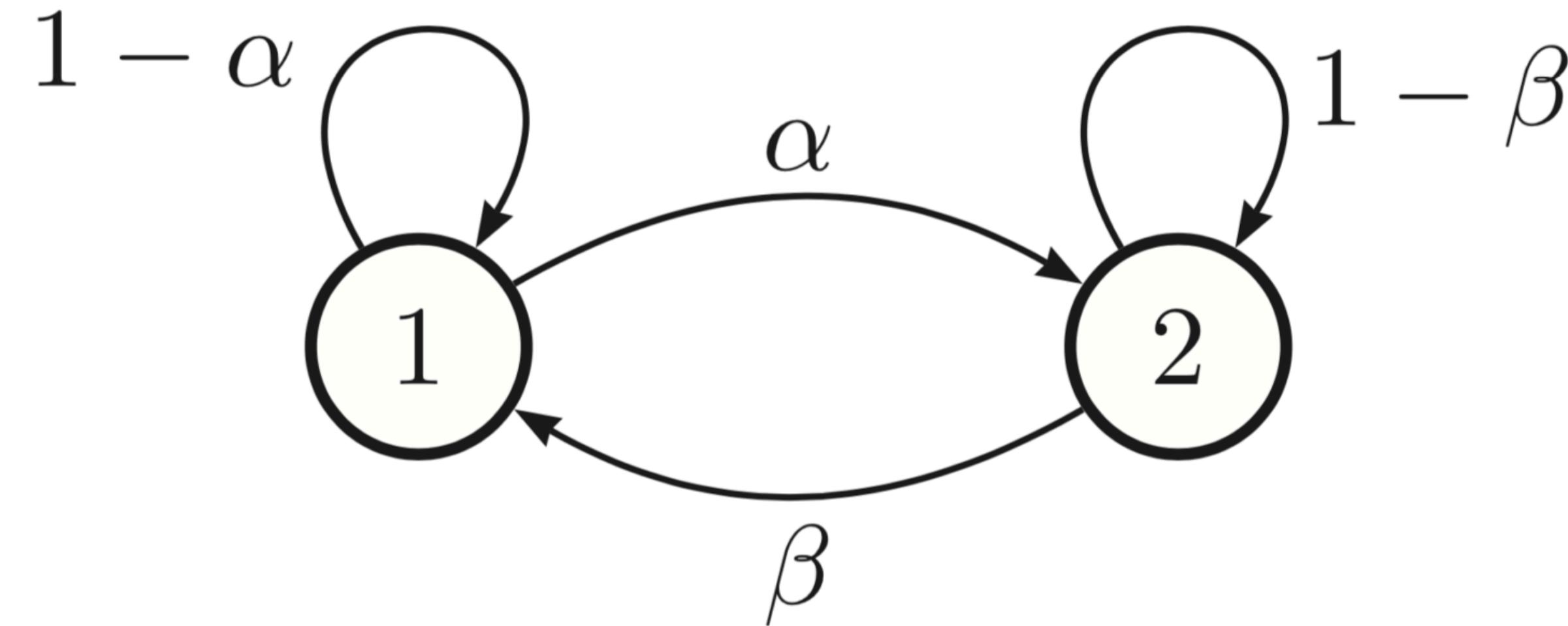
Recap: Markov chains

- **Markov Chains are probabilistic models for series that rely on Markov assumption**
- **Markov Assumption** for the first order Markov Chain can be written as X_t captures all relevant information to describe X_{t+1}

$$X_{t+1} | X_t \perp\!\!\!\perp X_{t-1}, X_{t-2}, \dots, X_1$$

- **in other words:**
- $P(X_{t+1} | X_t, X_{t-1}, \dots, X_1) = P(X_{t+1} | X_t)$

Transition Matrix



$$A = \begin{bmatrix} 1 - \alpha & \alpha \\ \beta & 1 - \beta \end{bmatrix}$$

- Quick check: The transition matrix after n steps looks like = ...

Markov chains. Example. Demo

- <http://setosa.io/blog/2014/07/26/markov-chains/>

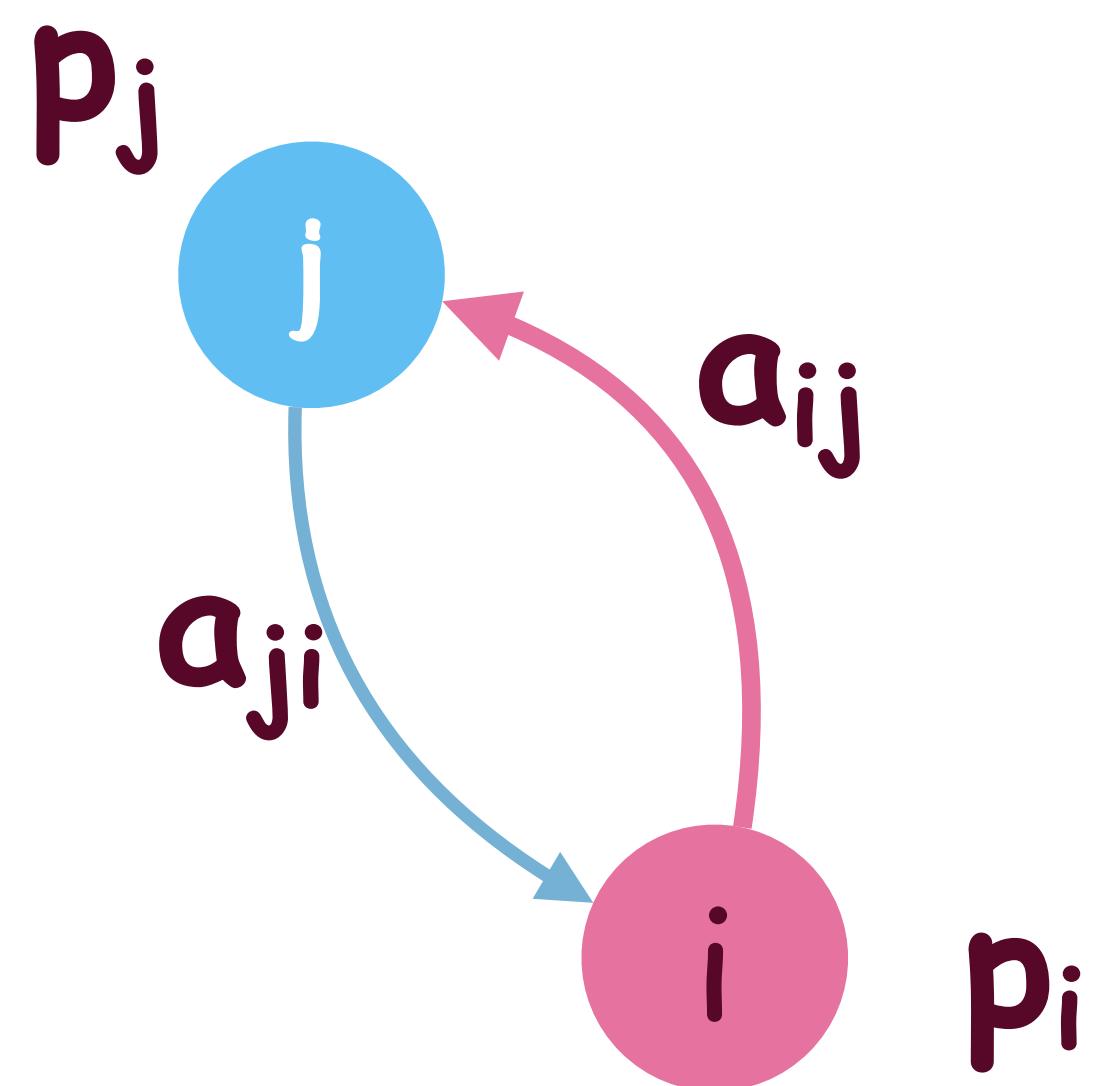
Stationary Distribution

- MCMC methods rely on the existence of stationary distribution for a Markov Chain.
- $p^\top = p^\top A$
 - p is a stationary distribution
 - A is a transition matrix
- we can run a Markov Chain for many steps and it may converge to stationary distr.
 - $p_0^\top A^\infty = p^\top$
 - p_0 is some random initialization

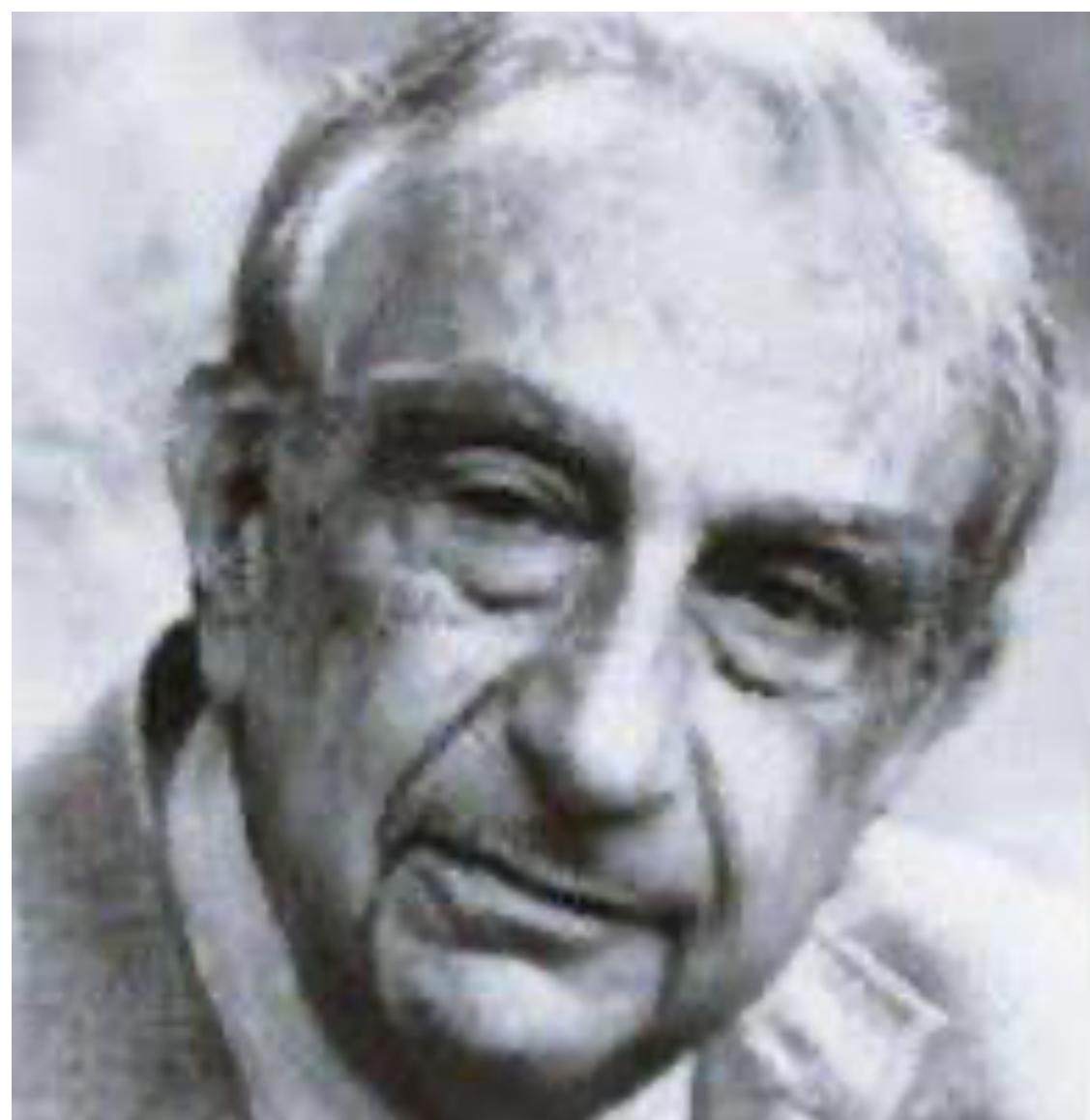
Conditions for a stationary distribution

- **necessary and sufficient conditions for stationary distribution to exist:**

- Markov chain is irreducible (it is possible to achieve any state from any state)
- a detailed balance condition: for all i,j :
- $\pi_i A_{ij} = \pi_j A_{ji}$



Metropolis-Hastings



Metropolis: Random walk

- **Random Walk:** $q(x^{(j)} | x) = q(x | x^{(j)})$

$$\rho(x^*, x^{(j)}) = \min \left(1, \frac{f(x^*)}{f(x^{(j)})} \right)$$

Metropolis-Hastings

- You have no access to CDF
- You cannot use Accept-Reject (too slow)
- But if you still can calculate p.d.f. up to a proportionality constant
 - Then you can use MCMC:
 - Metropolis Sampling
 - Metropolis-Hastings Sampling
 - Gibbs Sampling
 - Hamiltonian Monte Carlo

Metropolis-Hastings: Method

- Given some current value $x(j)$ sample the next value x^* using a proposal distribution $q(x)$:
 - sample x^* comes from $q(x | x^{(j)})$
- Calculate acceptance probability:
 - $$A(x^*, x^{(j)}) = \min \left(1, \frac{f(x^*)}{f(x^{(j)})} \frac{q(x^{(j)} | x^*)}{q(x^* | x^{(j)})} \right)$$
- Set $x(j+1) \leftarrow x^*$ with probability
$$A(x^*, x^{(j)})$$
- Otherwise, set $x(j+1) \leftarrow x(j)$

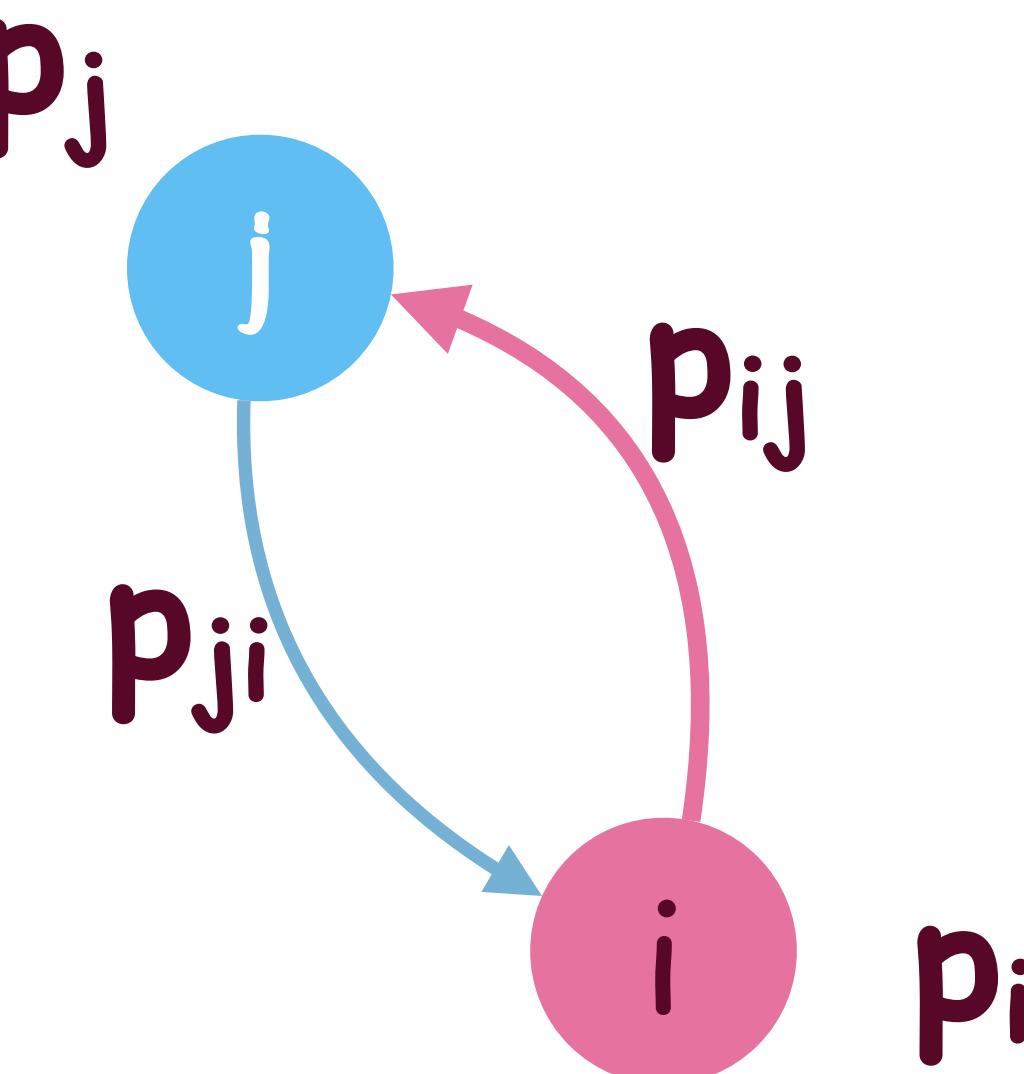
Metropolis-Hastings algorithm. discrete distribution

- The Metropolis-Hastings algorithm is a general method to design a Markov chain whose stationary distribution is a given target distribution p .
 - Start with a connected undirected graph G on the set of states.
- In general, let r be the maximum degree of any vertex of G .

Metropolis-Hastings algorithm. discrete distribution

- The transitions of the Markov chain are defined as follows:
- At state i select neighbour j with probability $\frac{1}{r}$.
- Since the degree of i may be less than r , with some probability no edge (transition) is selected and the walk remains at i .
- If a neighbour j is selected and $p_j \geq p_i$, go to j .
- If $p_j < p_i$, then

go to j with probability $\frac{p_j}{p_i}$ and
stay at i with probability $(1 - \frac{p_j}{p_i})$



Metropolis-Hastings algorithm

- The transitions of the Markov chain are defined as follows:
- At state i select neighbour j with probability $\frac{1}{r}$.
- Since the degree of i may be less than r , with some probability no edge (transition) is selected and the walk remains at i .
- If a neighbour j is selected and $p_j \geq p_i$, go to j .
- If $p_j < p_i$, then

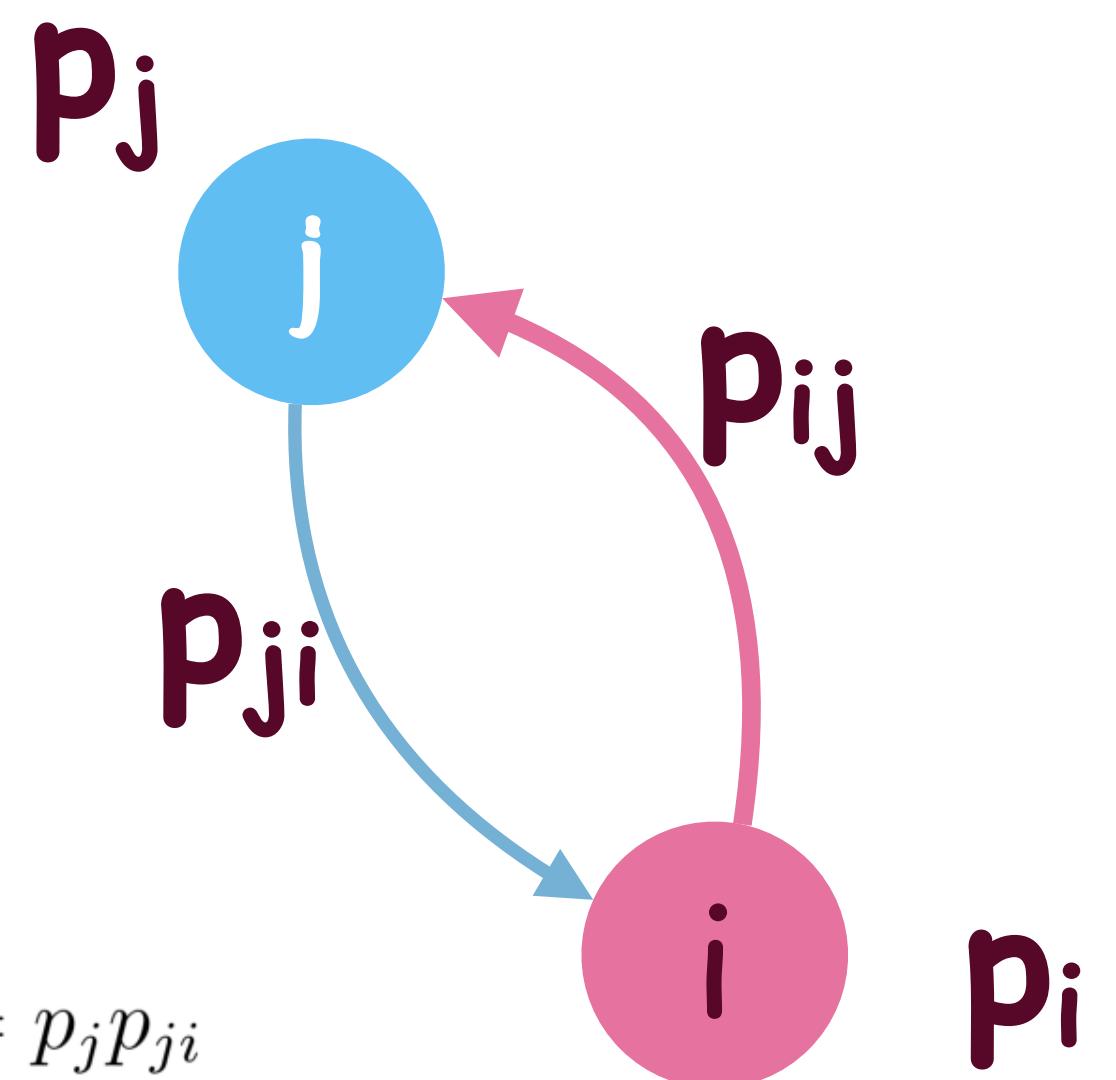
go to j with probability $\frac{p_j}{p_i}$ and

stay at i with probability $(1 - \frac{p_j}{p_i})$

$$p_i p_{ij} = \frac{p_i}{r} \min\left(1, \frac{p_j}{p_i}\right) = \frac{1}{r} \min(p_i, p_j) = \frac{p_j}{r} \min\left(1, \frac{p_i}{p_j}\right) = p_j p_{ji}$$

$$p_{ij} = \frac{1}{r} \min\left(1, \frac{p_j}{p_i}\right)$$

$$p_{ii} = 1 - \sum_{j \neq i} p_{ij}$$

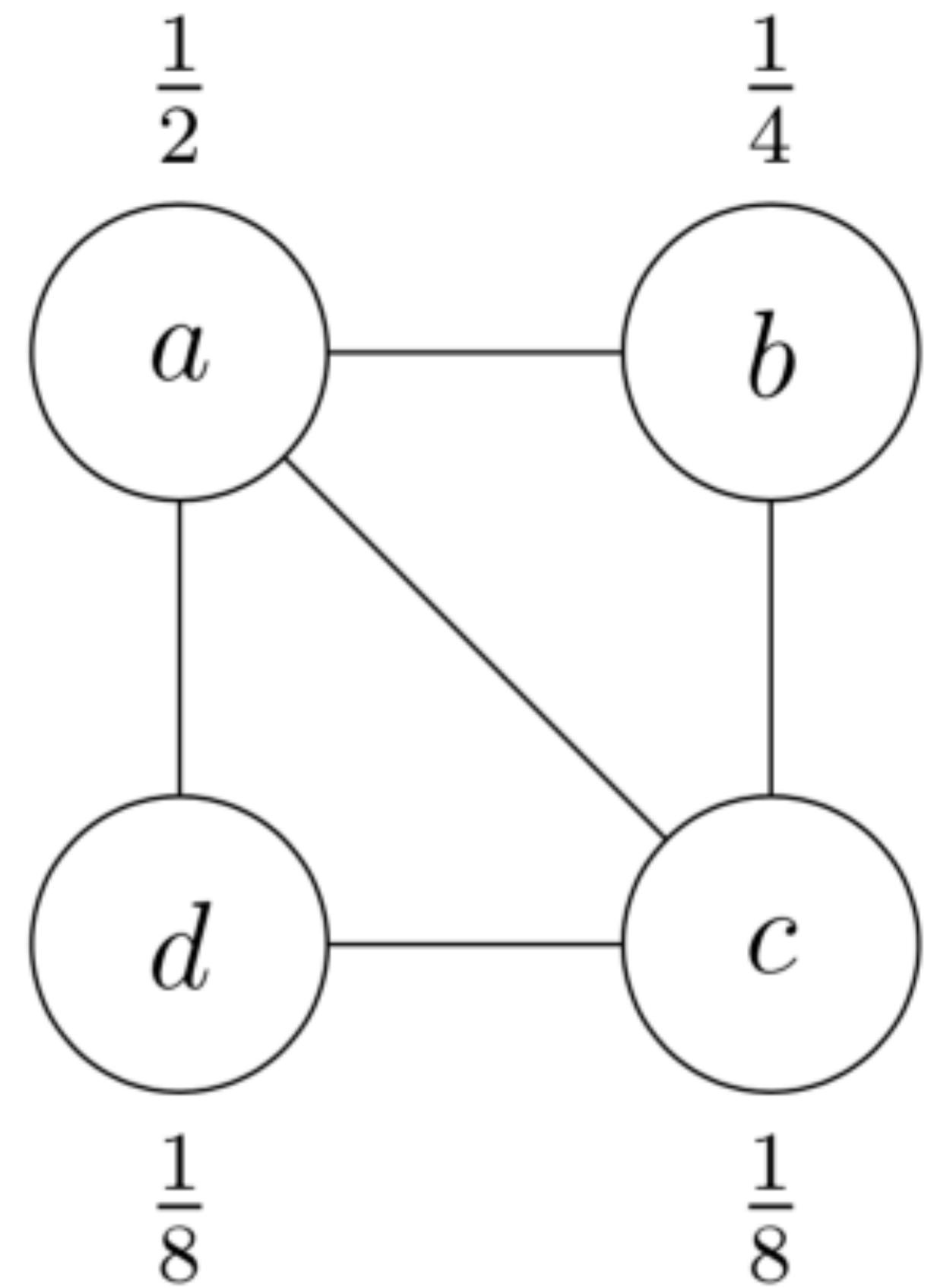


Example

- **Goals:**

- Construct the transition matrix using M-H
- Check the detailed balance property

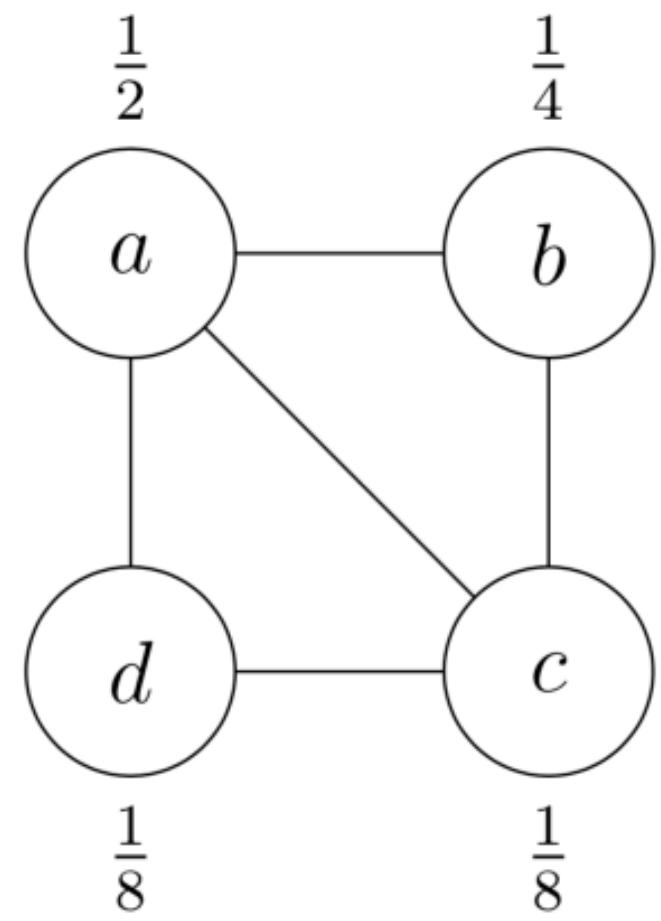
$$\begin{aligned} p(a) &= \frac{1}{2} \\ p(b) &= \frac{1}{4} \\ p(c) &= \frac{1}{8} \\ p(d) &= \frac{1}{8} \end{aligned}$$



Example

• **r=3**

$$\begin{aligned} p(a) &= \frac{1}{2} \\ p(b) &= \frac{1}{4} \\ p(c) &= \frac{1}{8} \\ p(d) &= \frac{1}{8} \end{aligned}$$



$a \rightarrow b$	$\frac{1}{3} \frac{1}{4} \frac{2}{1} = \frac{1}{6}$	$c \rightarrow a$	$\frac{1}{3}$
$a \rightarrow c$	$\frac{1}{3} \frac{1}{8} \frac{2}{1} = \frac{1}{12}$	$c \rightarrow b$	$\frac{1}{3}$
$a \rightarrow d$	$\frac{1}{3} \frac{1}{8} \frac{2}{1} = \frac{1}{12}$	$c \rightarrow d$	$\frac{1}{3}$
$a \rightarrow a$	$1 - \frac{1}{6} - \frac{1}{12} - \frac{1}{12} = \frac{2}{3}$	$c \rightarrow c$	$1 - \frac{1}{3} - \frac{1}{3} - \frac{1}{3} = 0$
$b \rightarrow a$	$\frac{1}{3}$	$d \rightarrow a$	$\frac{1}{3}$
$b \rightarrow c$	$\frac{1}{3} \frac{1}{8} \frac{4}{1} = \frac{1}{6}$	$d \rightarrow c$	$\frac{1}{3}$
$b \rightarrow b$	$1 - \frac{1}{3} - \frac{1}{6} = \frac{1}{2}$	$d \rightarrow d$	$1 - \frac{1}{3} - \frac{1}{3} = \frac{1}{3}$

$$\begin{aligned} p(a) &= p(a)p(a \rightarrow a) + p(b)p(b \rightarrow a) + p(c)p(c \rightarrow a) + p(d)p(d \rightarrow a) \\ &= \frac{1}{2} \frac{2}{3} + \frac{1}{4} \frac{1}{3} + \frac{1}{8} \frac{1}{3} + \frac{1}{8} \frac{1}{3} = \frac{1}{2} \end{aligned}$$

Check the property

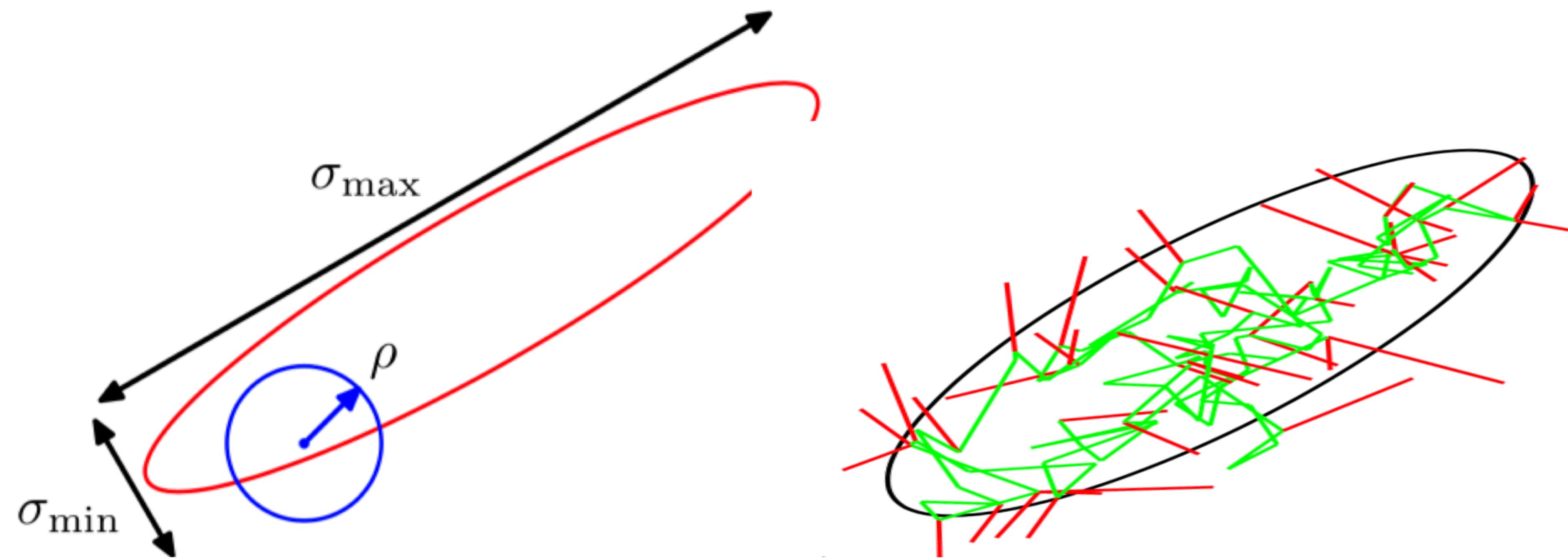
$$\begin{aligned} p(a) &= p(a)p(a \rightarrow a) + p(b)p(b \rightarrow a) + p(c)p(c \rightarrow a) + p(d)p(d \rightarrow a) \\ &= \frac{1}{2} \cdot \frac{2}{3} + \frac{1}{4} \cdot \frac{1}{3} + \frac{1}{8} \cdot \frac{1}{3} + \frac{1}{8} \cdot \frac{1}{3} = \frac{1}{2} \end{aligned}$$

$$\begin{aligned} p(b) &= p(a)p(a \rightarrow b) + p(b)p(b \rightarrow b) + p(c)p(c \rightarrow b) \\ &= \frac{1}{2} \cdot \frac{1}{6} + \frac{1}{4} \cdot \frac{1}{2} + \frac{1}{8} \cdot \frac{1}{3} = \frac{1}{4} \end{aligned}$$

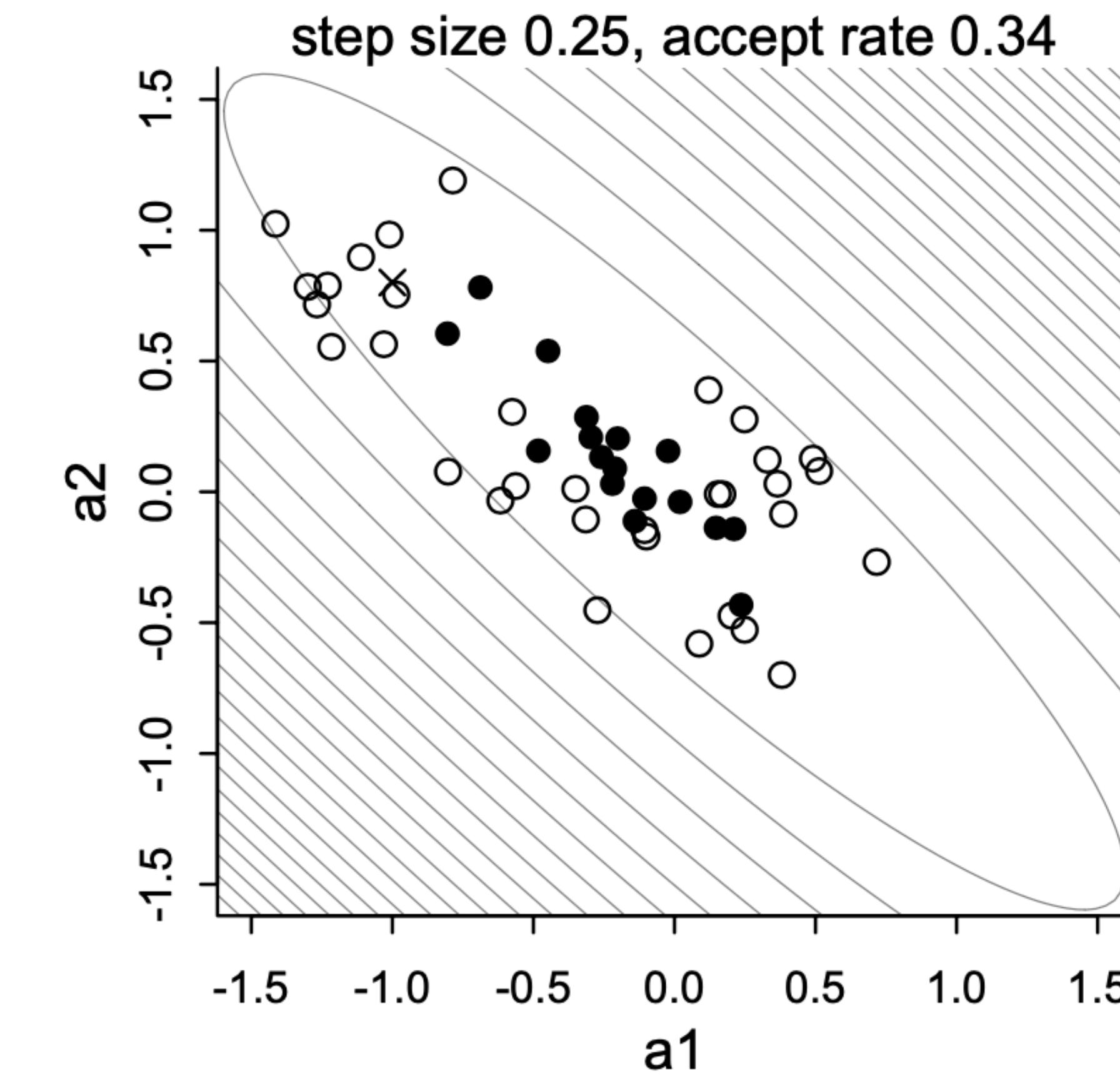
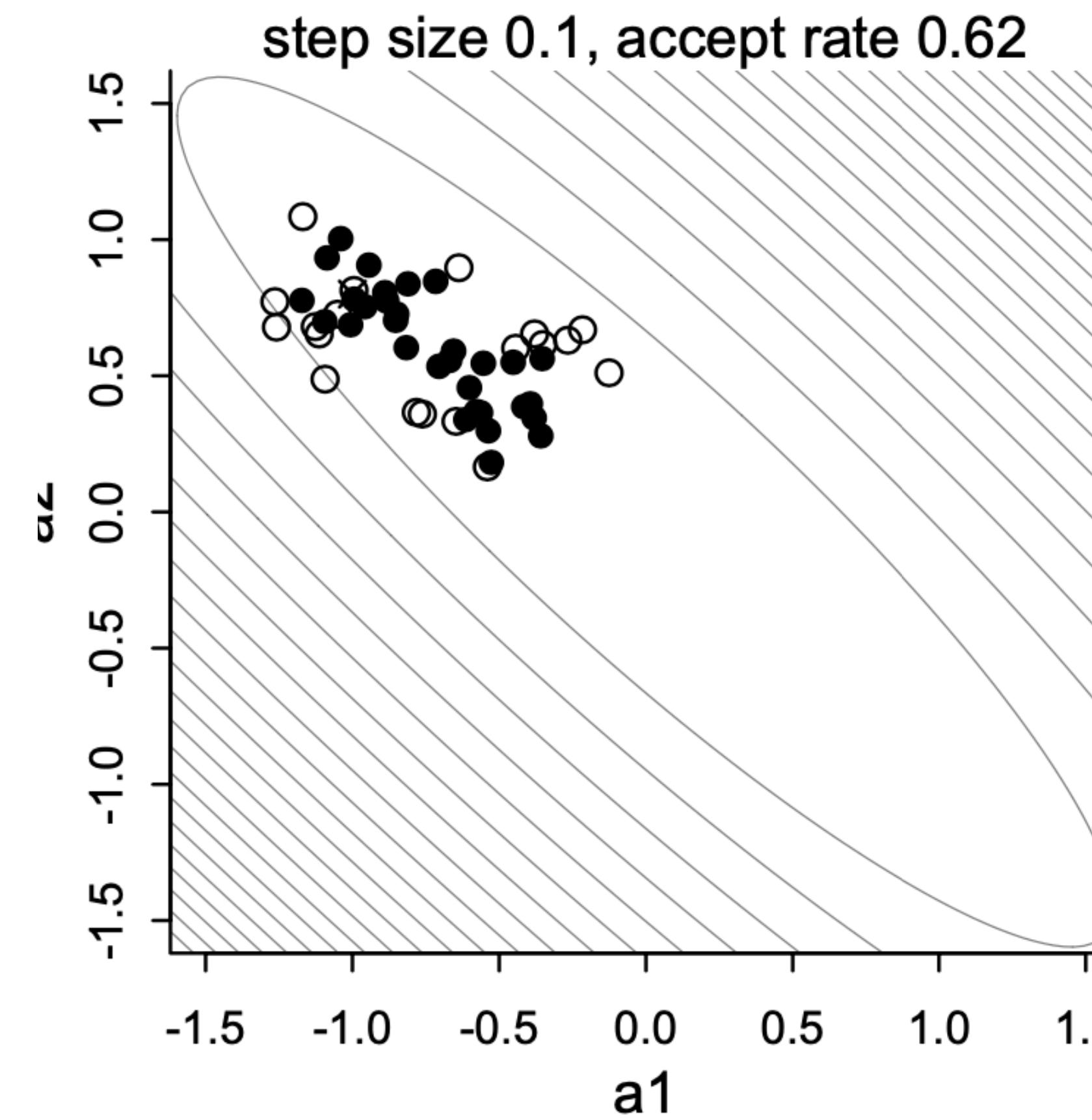
$$\begin{aligned} p(c) &= p(a)p(a \rightarrow c) + p(b)p(b \rightarrow c) + p(c)p(c \rightarrow c) + p(d)p(d \rightarrow c) \\ &= \frac{1}{2} \cdot \frac{1}{12} + \frac{1}{4} \cdot \frac{1}{6} + \frac{1}{8} \cdot 0 + \frac{1}{8} \cdot \frac{1}{3} = \frac{1}{8} \end{aligned}$$

$$\begin{aligned} p(d) &= p(a)p(a \rightarrow d) + p(c)p(c \rightarrow d) + p(d)p(d \rightarrow d) \\ &= \frac{1}{2} \cdot \frac{1}{12} + \frac{1}{8} \cdot \frac{1}{3} + \frac{1}{8} \cdot \frac{1}{3} = \frac{1}{8} \end{aligned}$$

Visualization of MH: continuous case

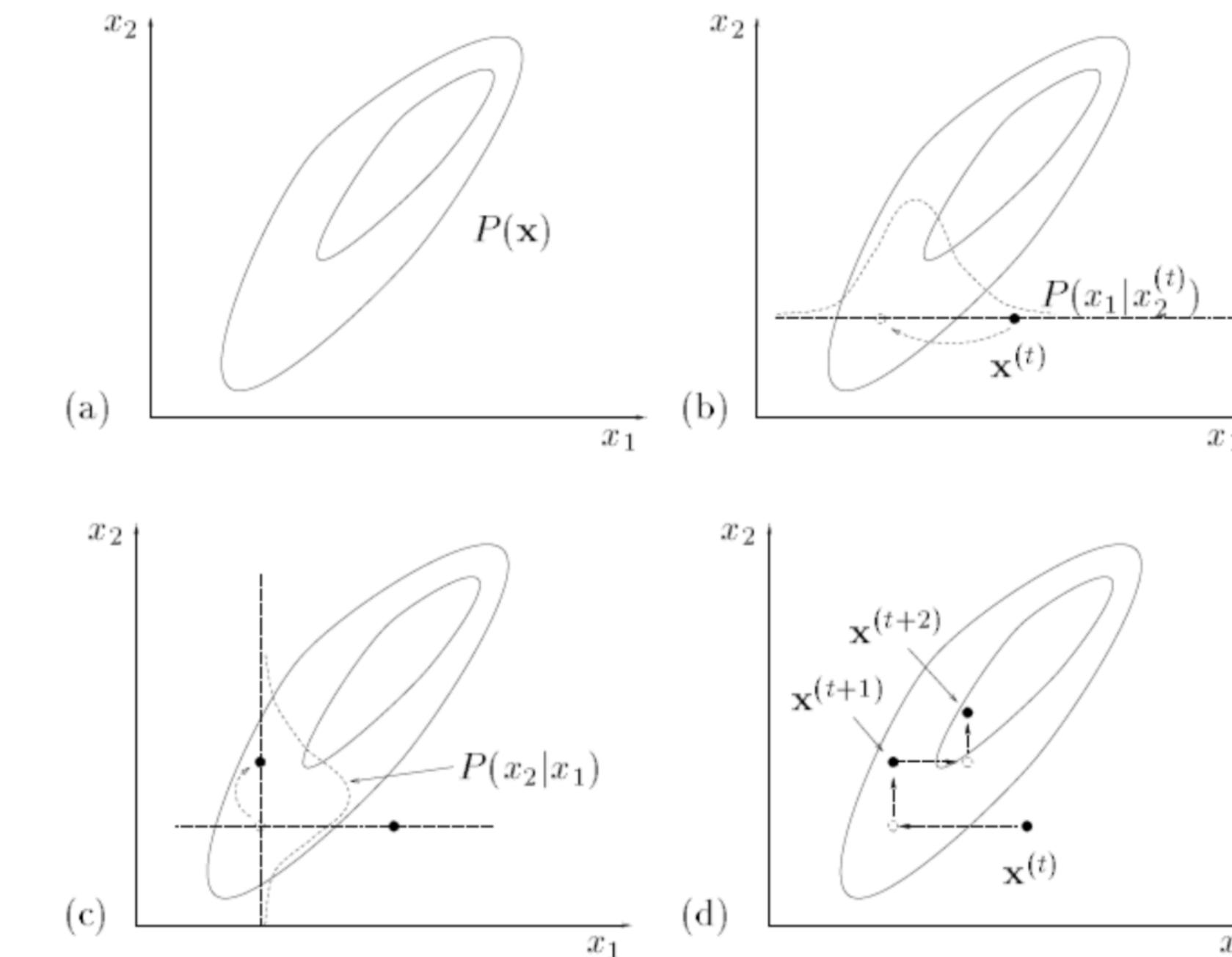
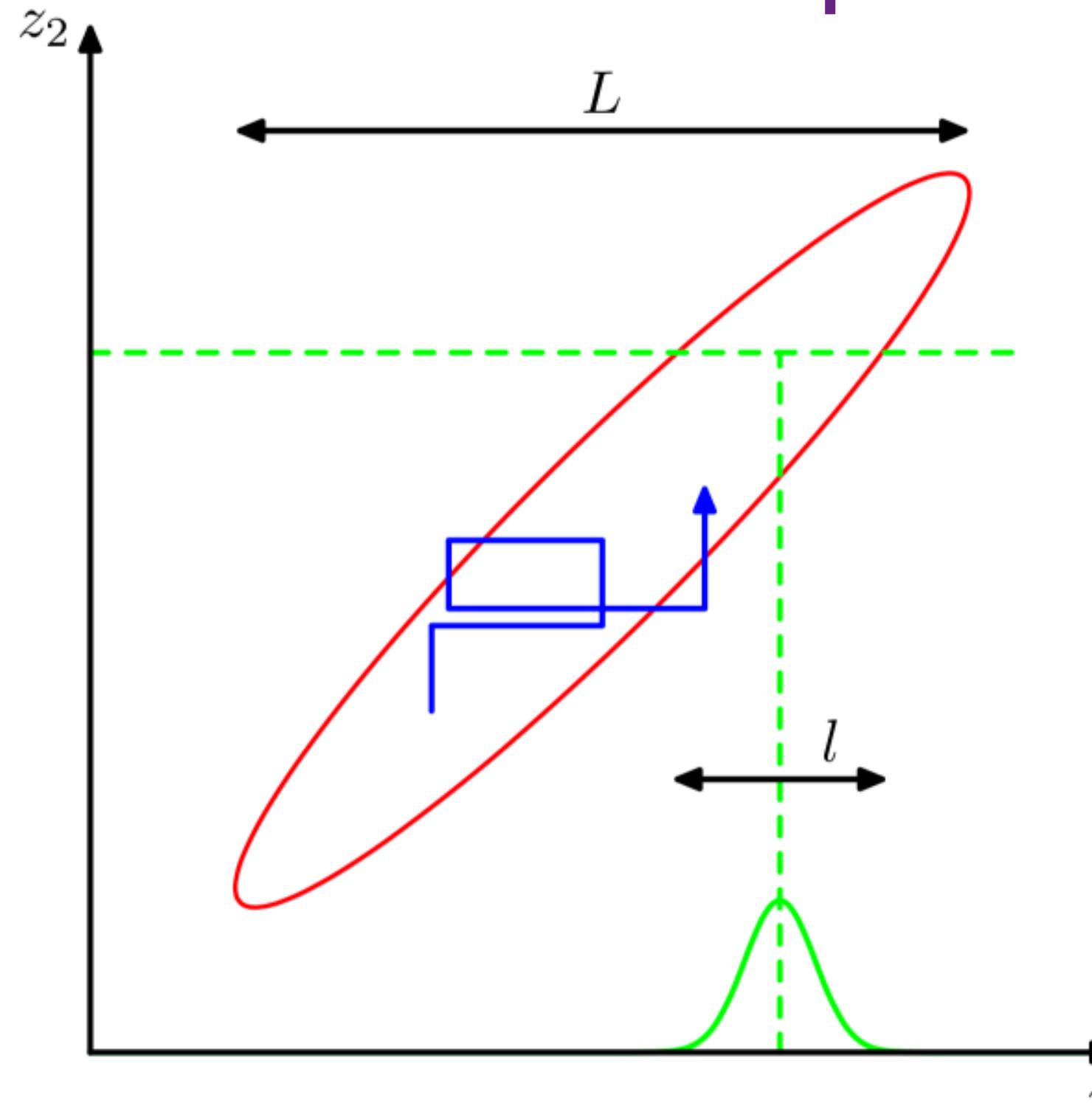


Metropolis chains under high correlation



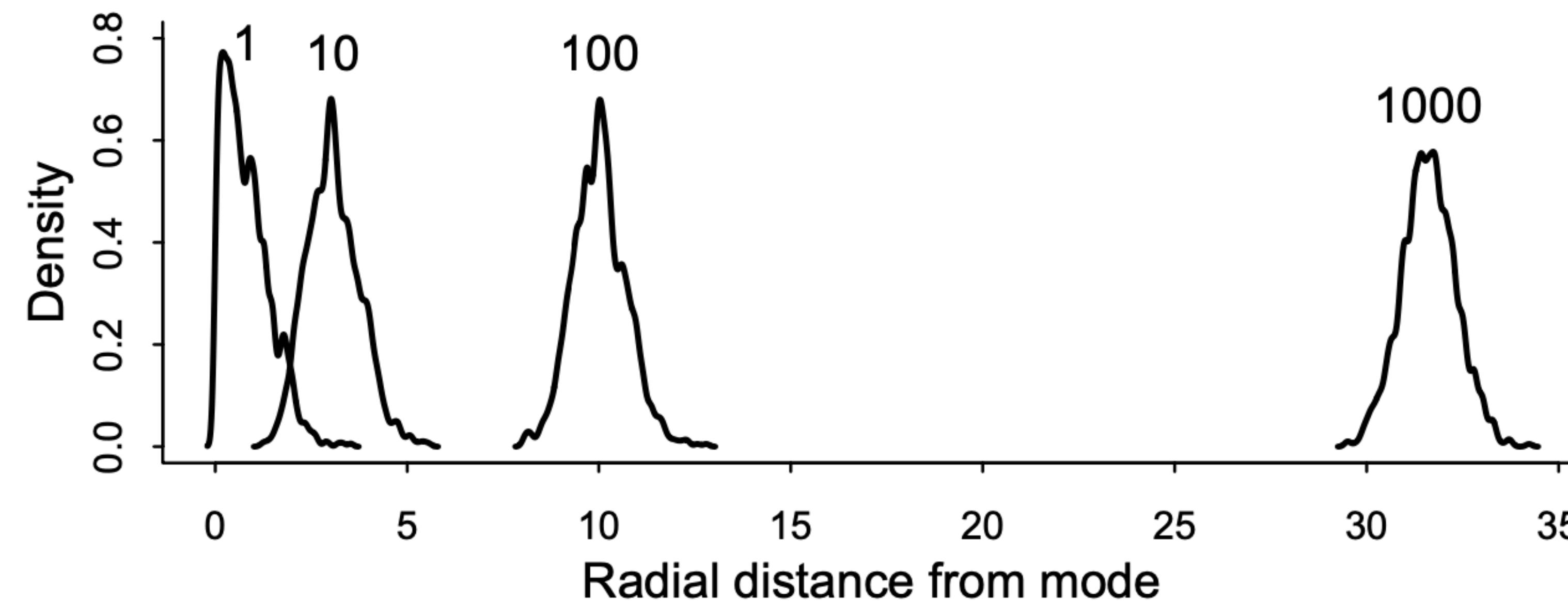
Gibbs Sampling

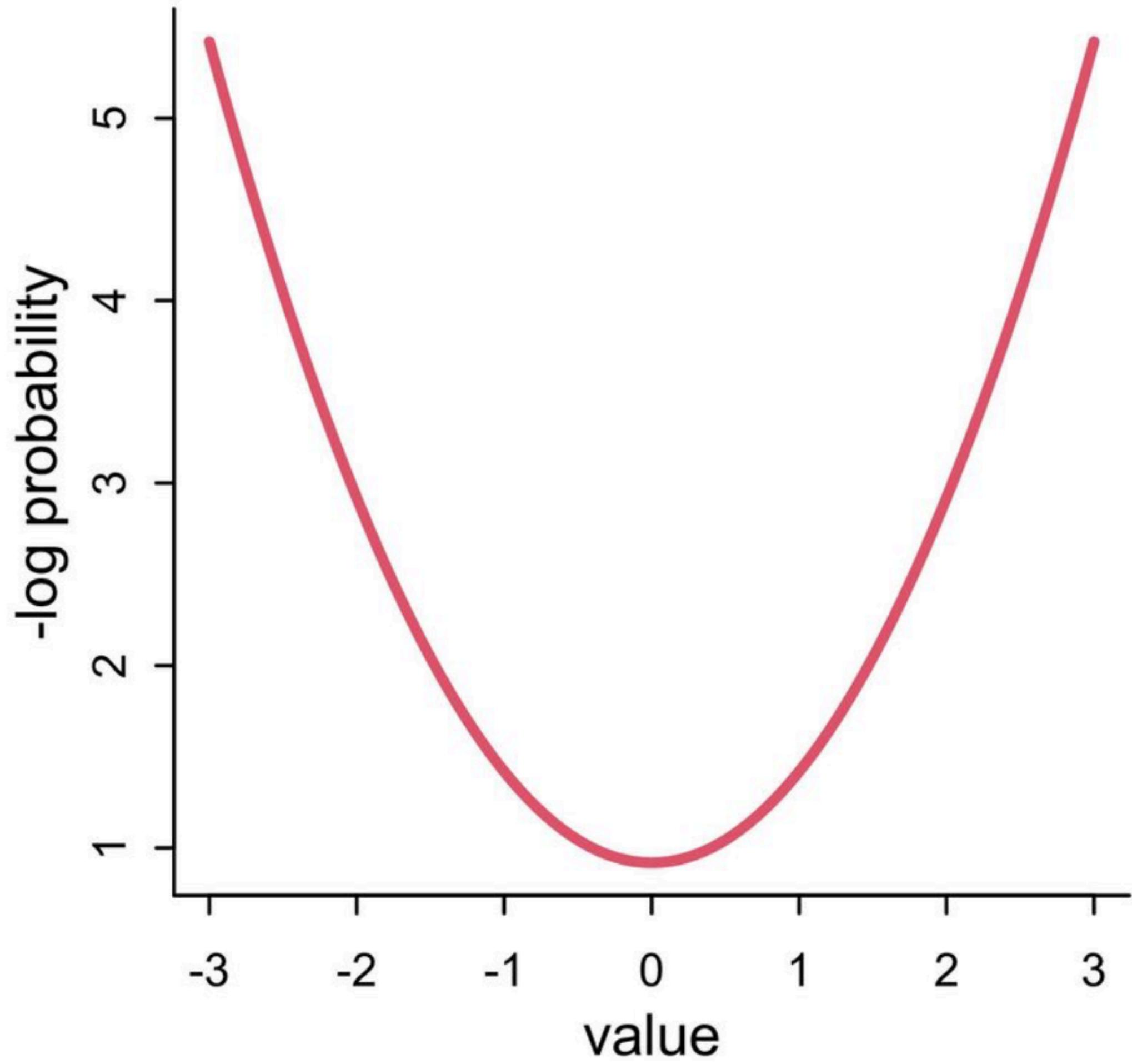
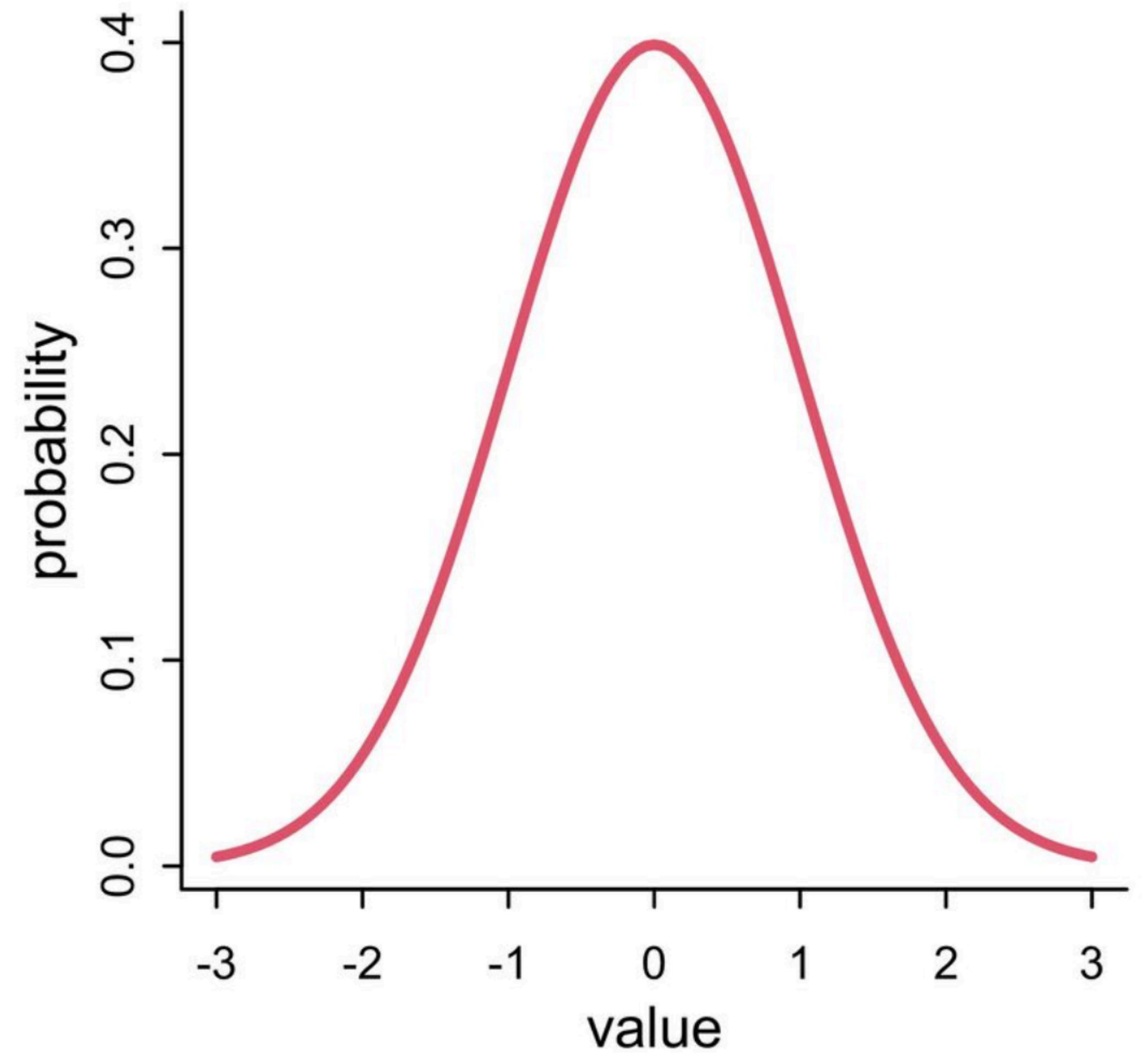
- The idea in Gibbs sampling is that, rather than probabilistically picking the next state all at once, you make a separate probabilistic choice for one of the d dimensions, where each choice depends on the other $d - 1$ dimensions.

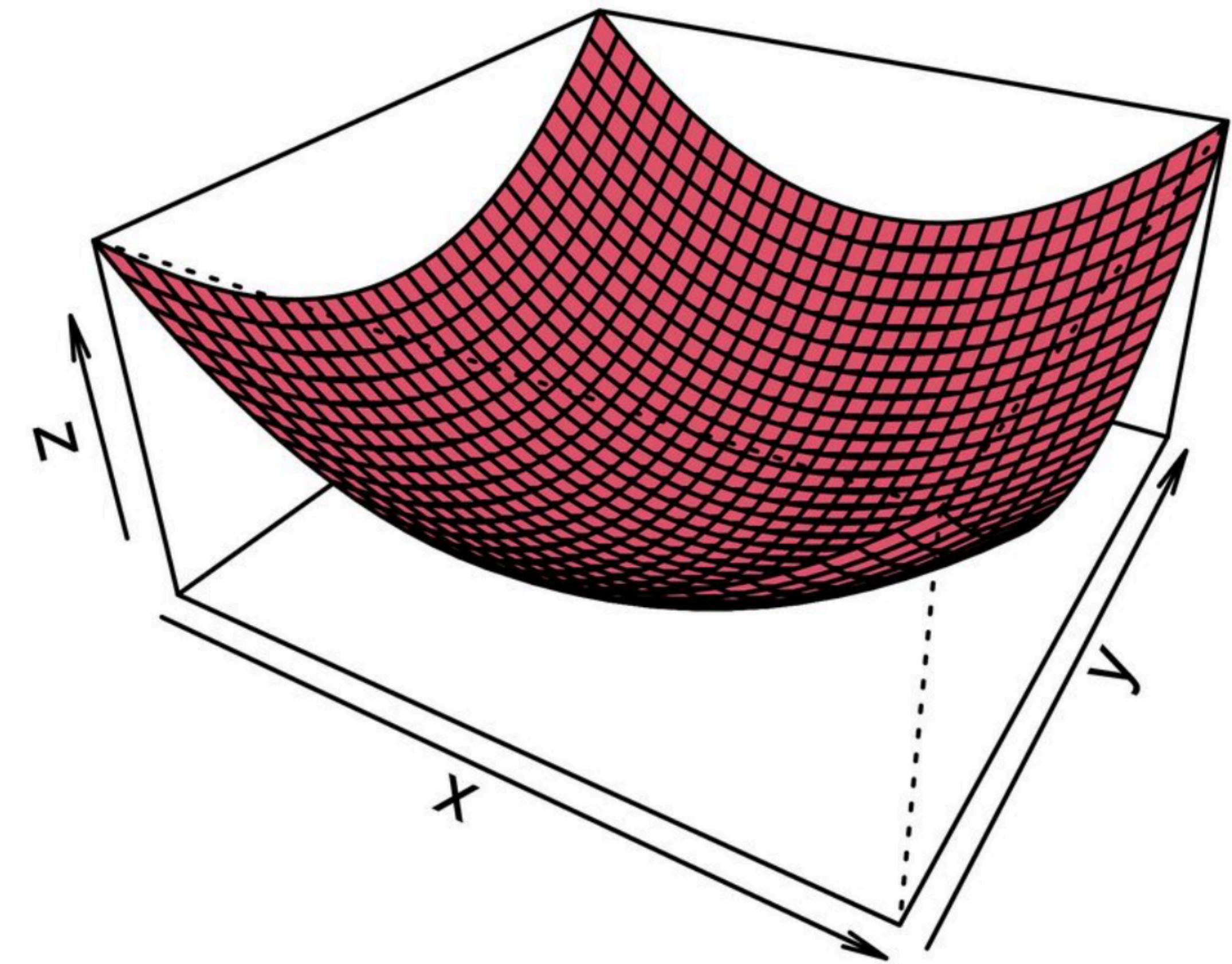
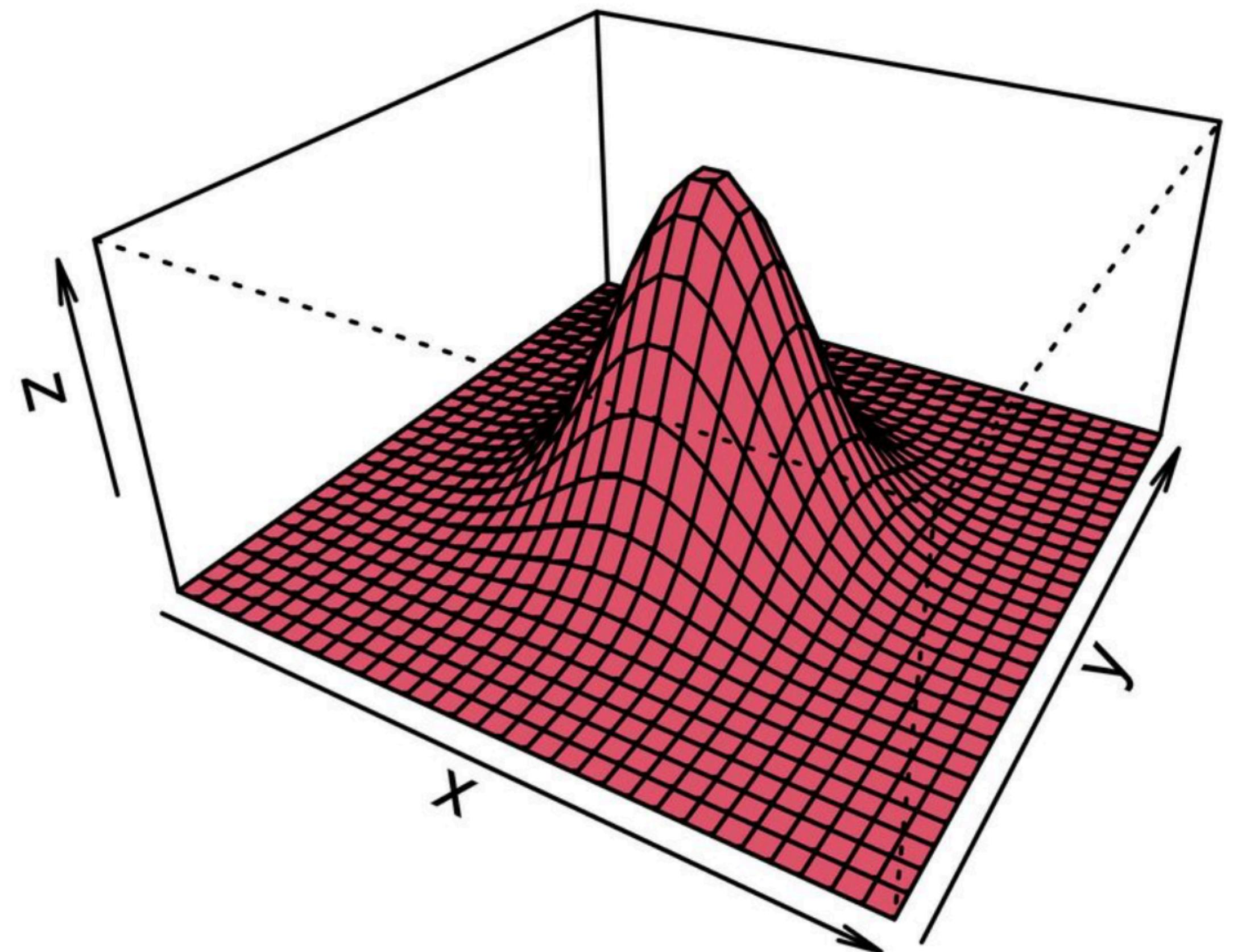


High-dimensional problems. Concentration of measure

- The most of the probability mass of a high-dimension distribution is always very far from the mode of the distribution.
- => The combination of parameter values that maximizes posterior probability, the mode, is not actually in a region of parameter values that are highly plausible



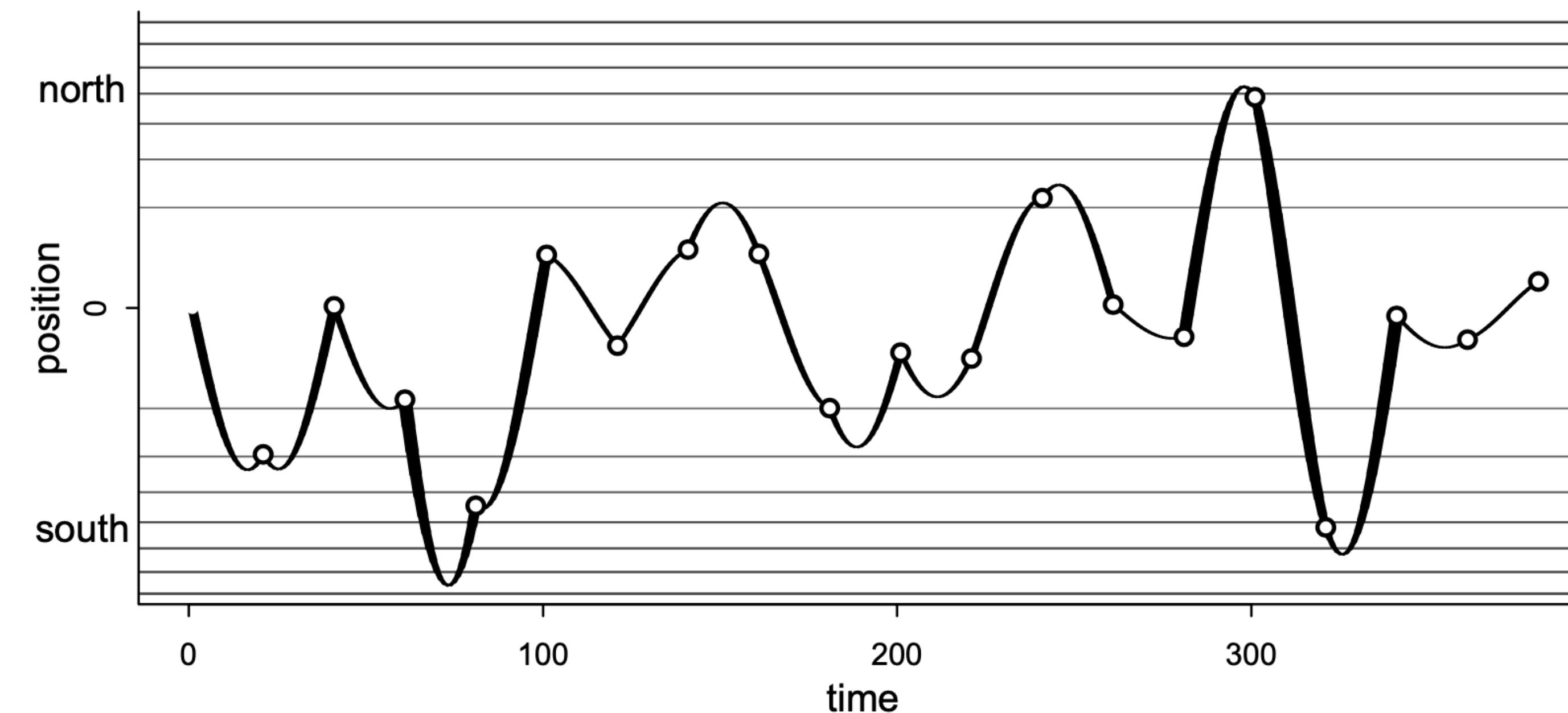




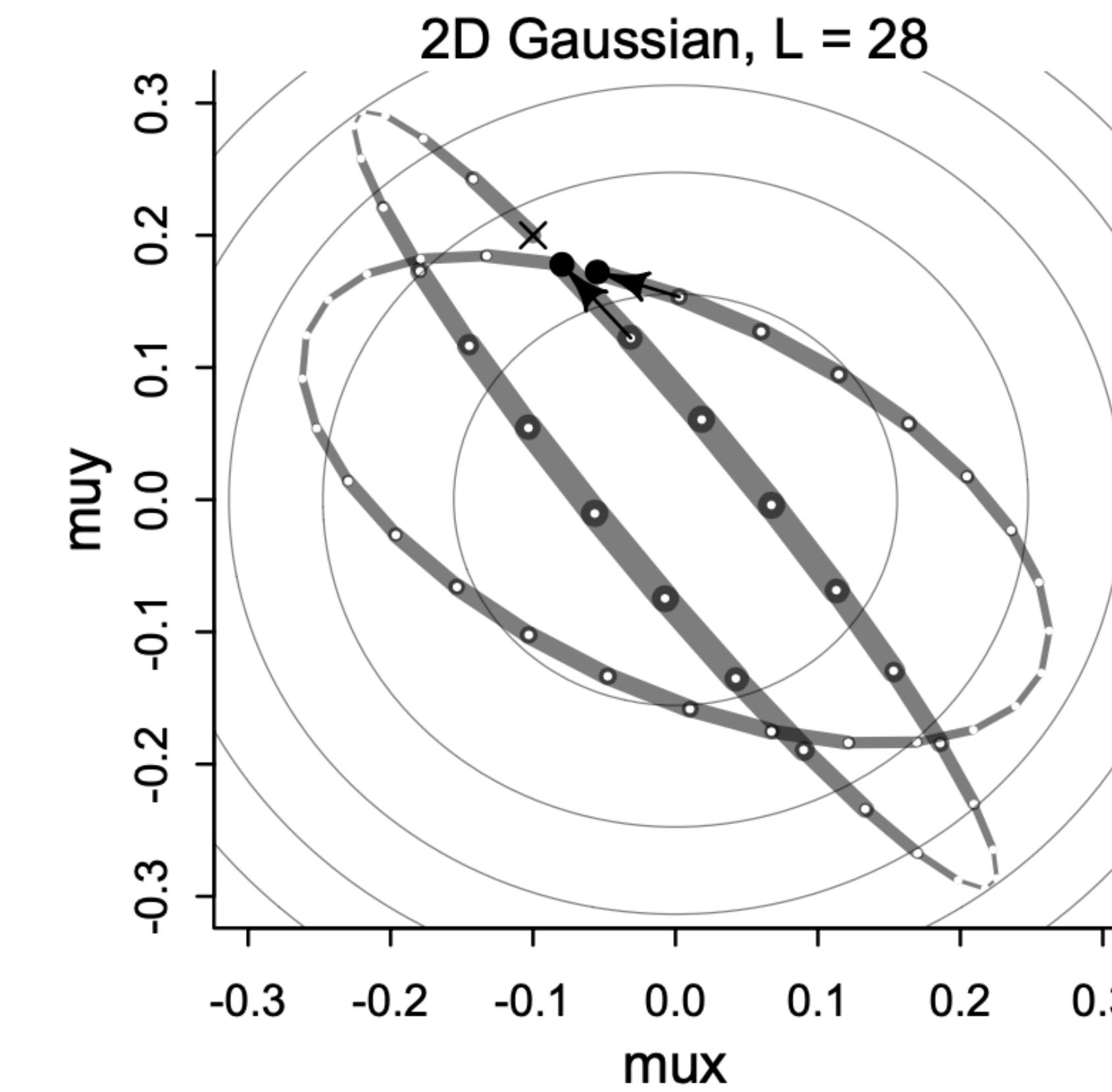
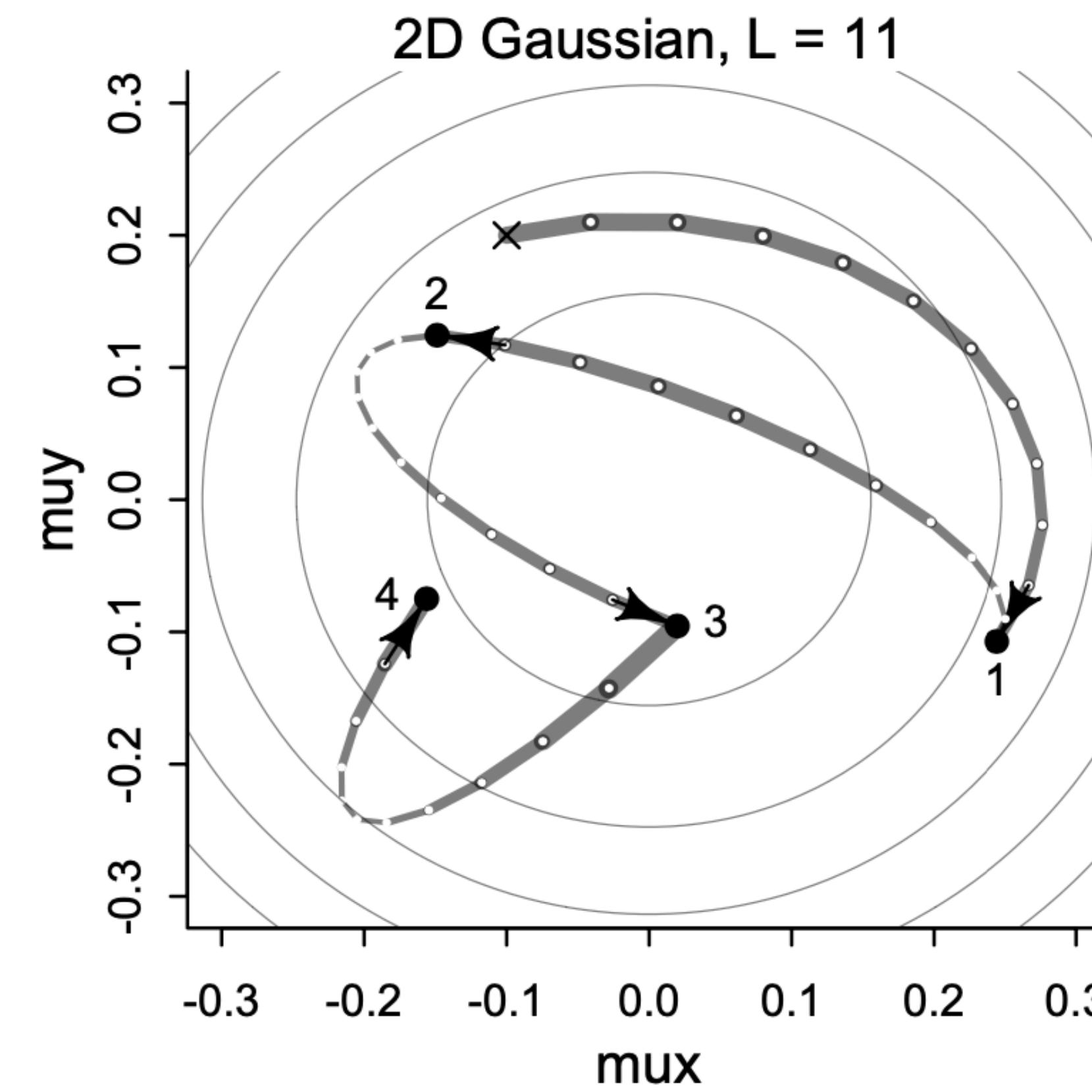


Hamiltonian Monte Carlo (Hybrid Monte Carlo, HMC)

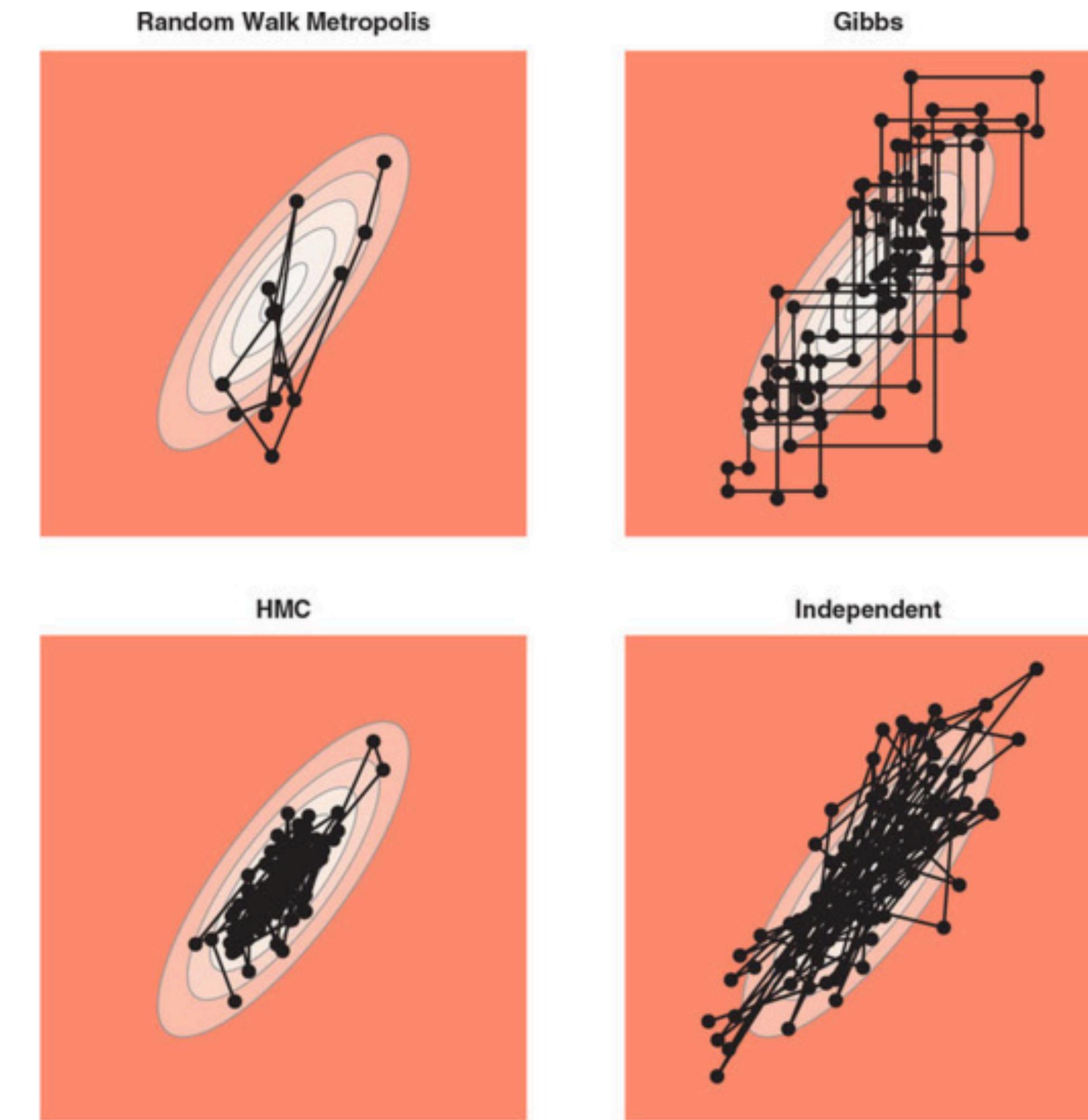
- HMC really does run a physics simulation, pretending the vector of parameters gives the position of a little frictionless particle.
- The log-posterior provides a surface for this particle to glide across.



U-turns



Different Sampling methods



- Lambert, Ben. *A Students Guide to Bayesian Statistics* SAGE Publications. 2018

Assignment: at home (optional)

- **Using Metropolis-Hastings**
 - Sample from $f(x)$:

$$f(x) \propto \exp(- (x - 3)^2 / 18) I(1 < x < 6)$$

- Let $q(x | x(j)) = N(3, 1)$; $x(0) = 2$
- Using samples, calculate the mean

References

- Probabilistic inference using Markov chain Monte Carlo methods, Radford M. Neal, Technical report: CRG-TR-93-1, Department of Computer Science, University of Toronto, 1993. <http://www.cs.toronto.edu/~radford/review.abstract.html>
- Various figures and more came from (see also references therein): Advances in Markov chain Monte Carlo methods. Iain Murray. 2007. <http://www.cs.toronto.edu/~murray/pub/07thesis/>
- Information theory, inference, and learning algorithms. David MacKay, 2003. <http://www.inference.phy.cam.ac.uk/mackay/itila/>
- Pattern recognition and machine learning. Christopher M. Bishop. 2006. <http://research.microsoft.com/~cmbishop/PRML/>
- Gibbs sampling for graphical models: <http://mathstat.helsinki.fi/openbugs/> <http://www-ice.iarc.fr/~martyn/software/jags/>
- Neural networks and other flexible models: <http://www.cs.utoronto.ca/~radford/fbm.software.html>
- CODA: <http://www-fis.iarc.fr/coda/>
-