

Statistical Techniques for Data Science & Robotics

Week 12



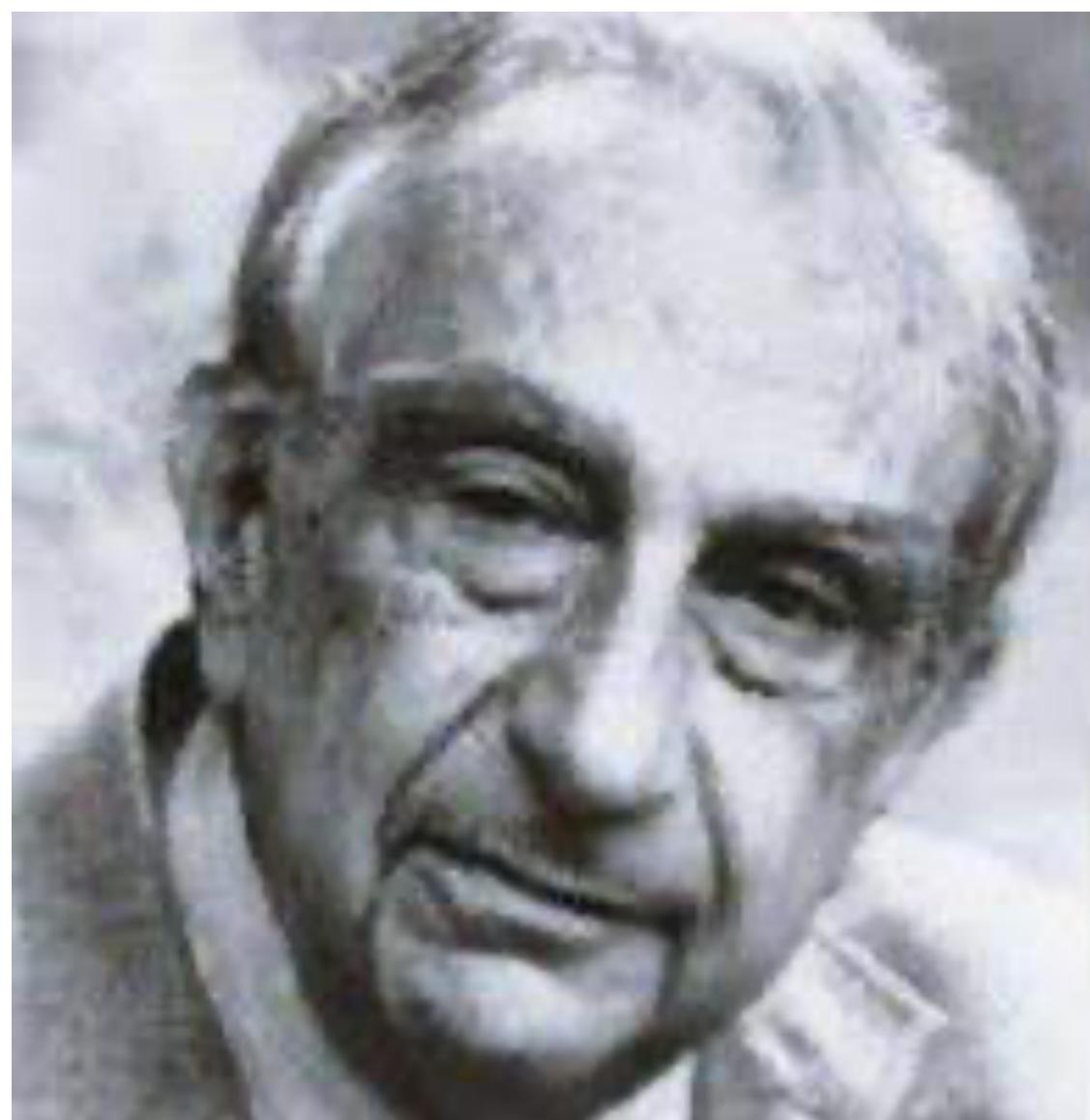
Quiz

- Q1: Describe **the detailed balance** property.

Objectives for today

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

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 $x(j+1) \leftarrow x(j)$

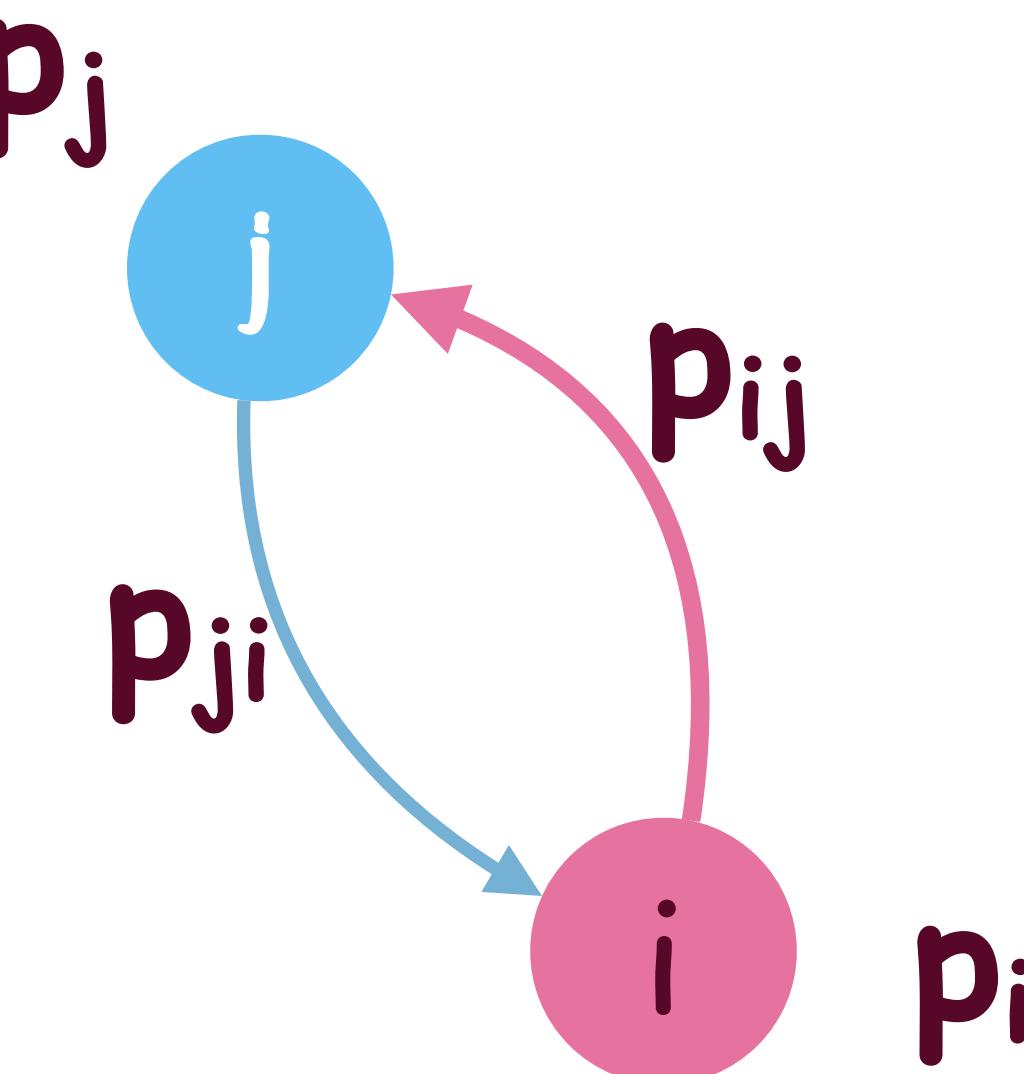
Metropolis-Hastings algorithm. discrete distribution

- 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

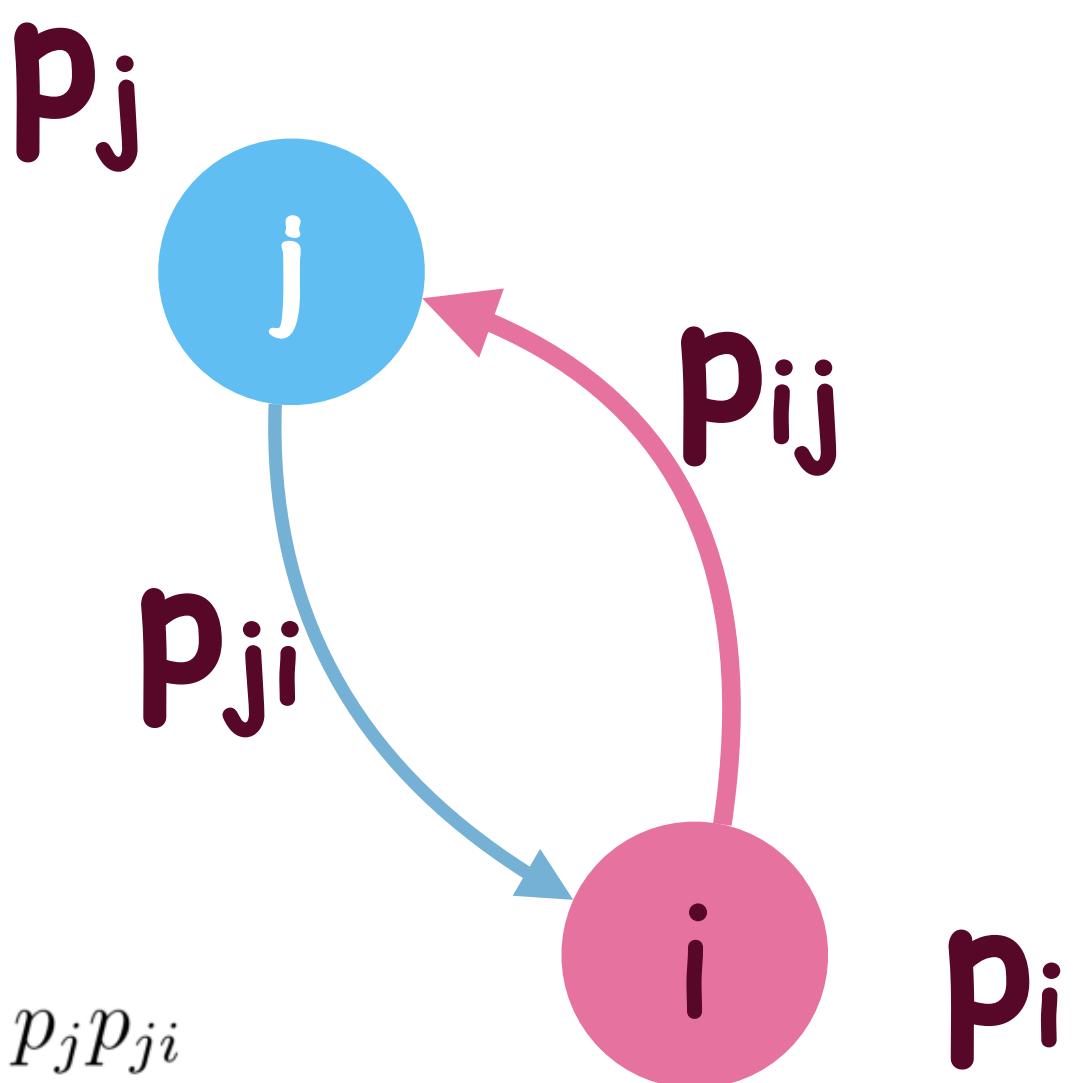
- The transitions of the Markov chain are defined as follows:
- At state i select neighbour j with probability $\frac{1}{r}$.
- Since the degree of i may be less than r , with some probability no edge (transition) is selected and the walk remains at i .
- If a neighbour j is selected and $p_j \geq p_i$,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

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

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



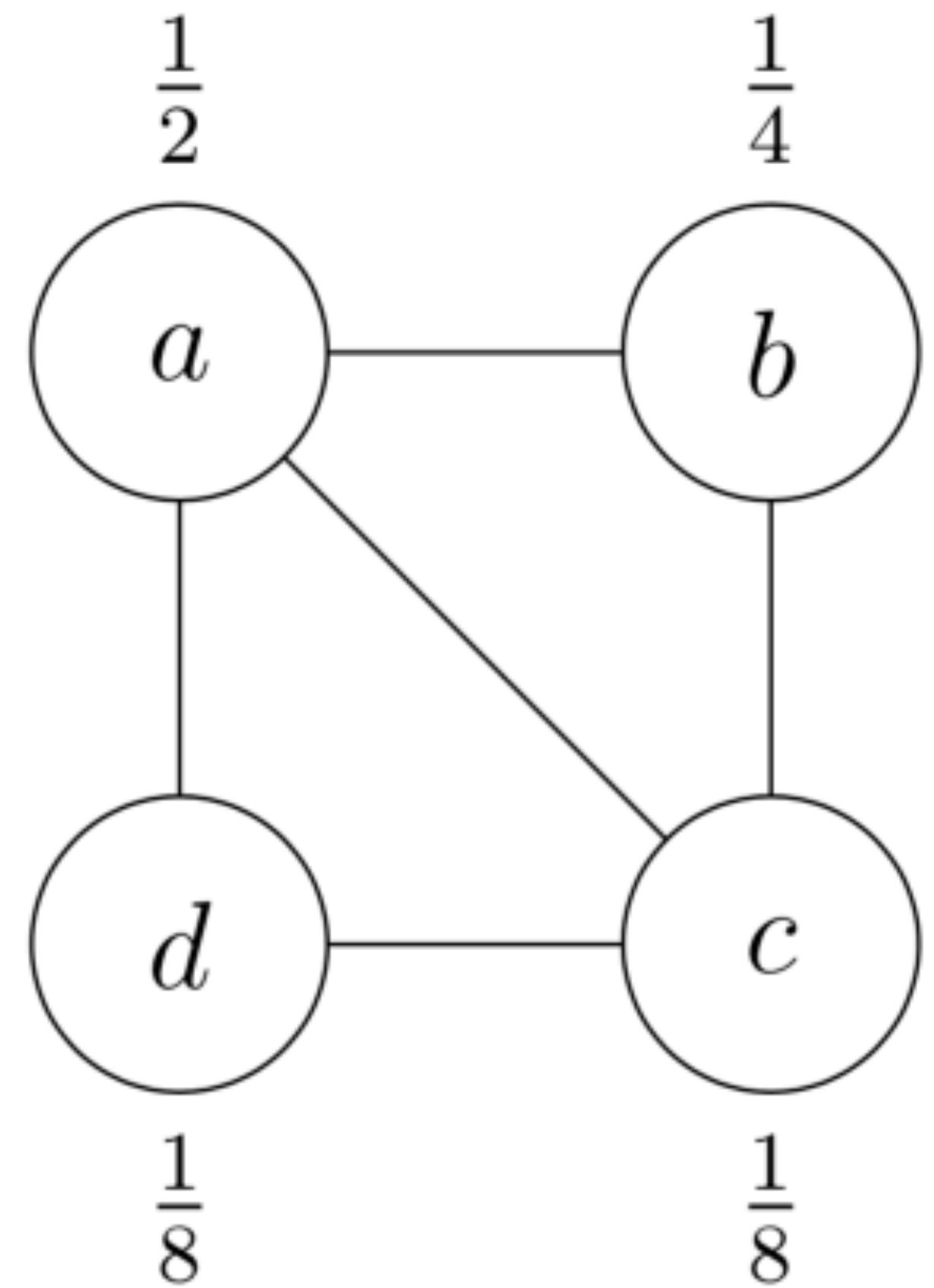
$$p_i p_{ij} = \frac{p_i}{r} \min \left(1, \frac{p_j}{p_i} \right) = \frac{1}{r} \min(p_i, p_j) = \frac{p_j}{r} \min \left(1, \frac{p_i}{p_j} \right) = p_j p_{ji}$$

Example

- **Goals:**

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

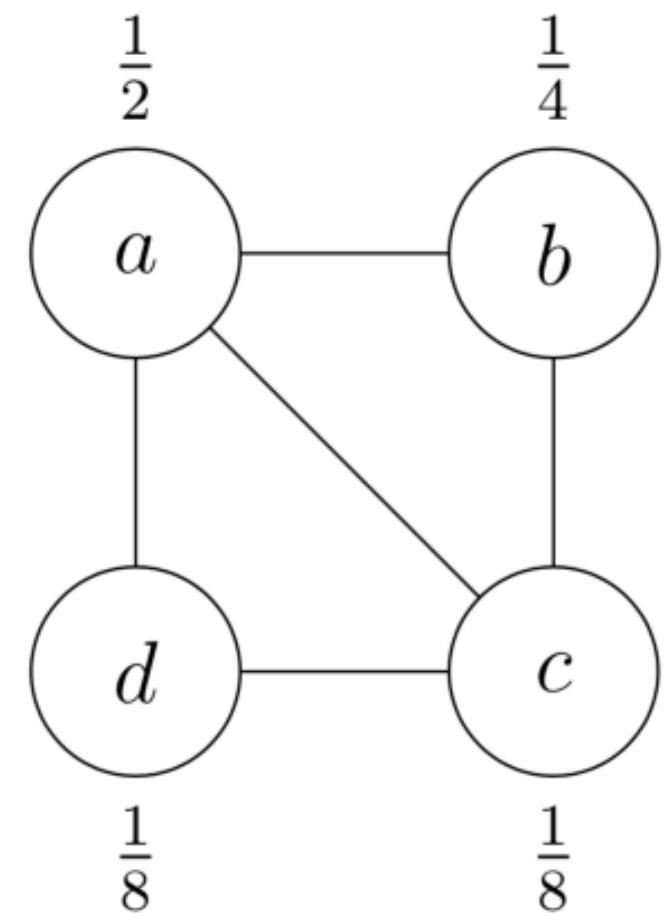
$$\begin{aligned} p(a) &= \frac{1}{2} \\ p(b) &= \frac{1}{4} \\ p(c) &= \frac{1}{8} \\ p(d) &= \frac{1}{8} \end{aligned}$$



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

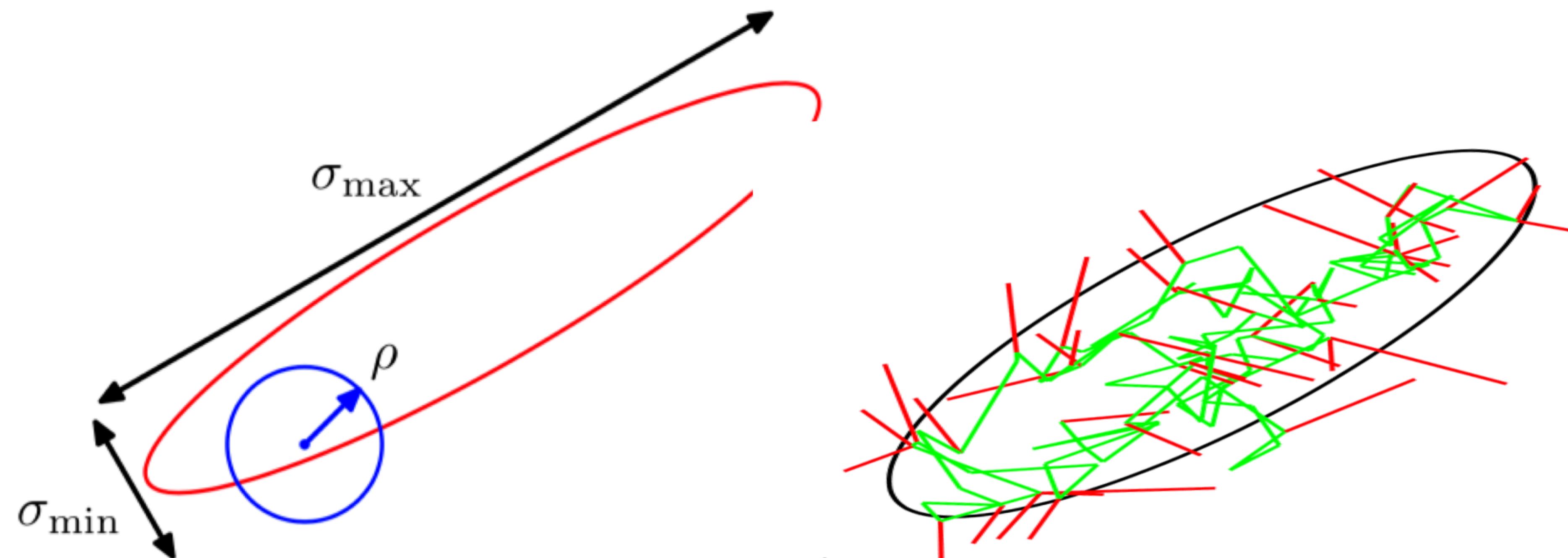
$$\begin{aligned} p(a) &= p(a)p(a \rightarrow a) + p(b)p(b \rightarrow a) + p(c)p(c \rightarrow a) + p(d)p(d \rightarrow a) \\ &= \frac{1}{2} \cdot \frac{2}{3} + \frac{1}{4} \cdot \frac{1}{3} + \frac{1}{8} \cdot \frac{1}{3} + \frac{1}{8} \cdot \frac{1}{3} = \frac{1}{2} \end{aligned}$$

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

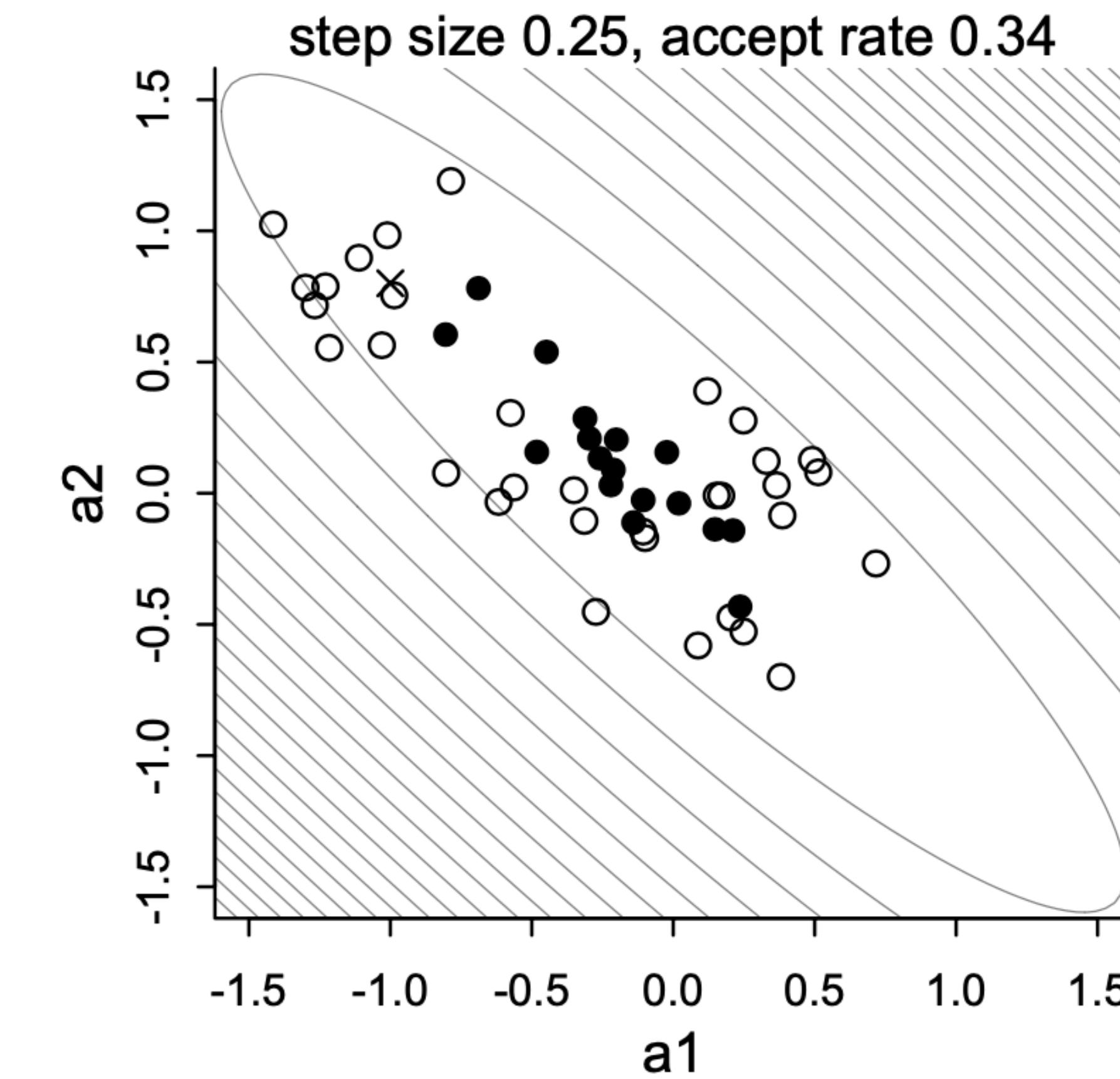
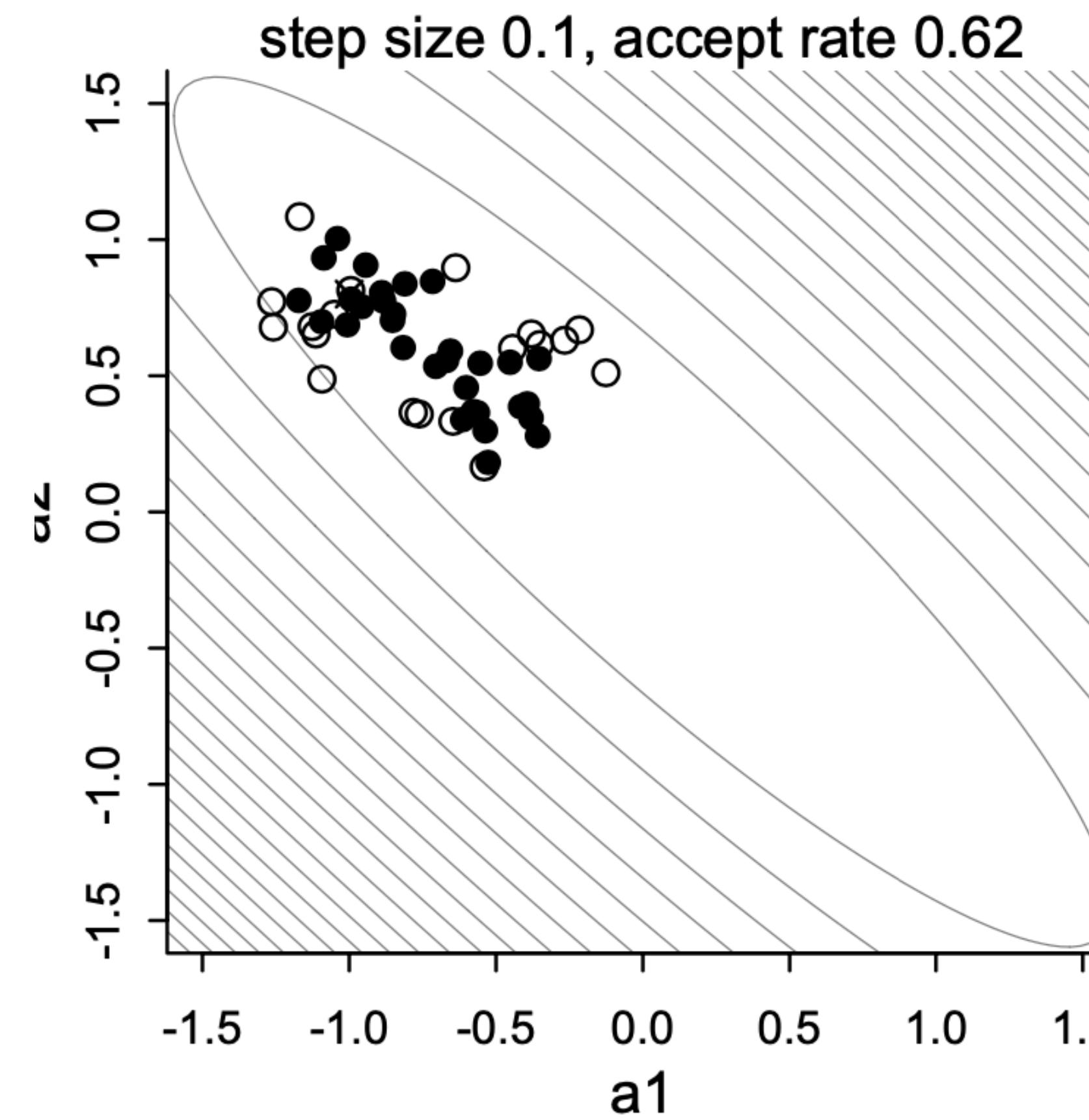
$$\begin{aligned} p(c) &= p(a)p(a \rightarrow c) + p(b)p(b \rightarrow c) + p(c)p(c \rightarrow c) + p(d)p(d \rightarrow c) \\ &= \frac{1}{2} \cdot \frac{1}{12} + \frac{1}{4} \cdot \frac{1}{6} + \frac{1}{8} \cdot 0 + \frac{1}{8} \cdot \frac{1}{3} = \frac{1}{8} \end{aligned}$$

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

Visualization of MH: continuous case

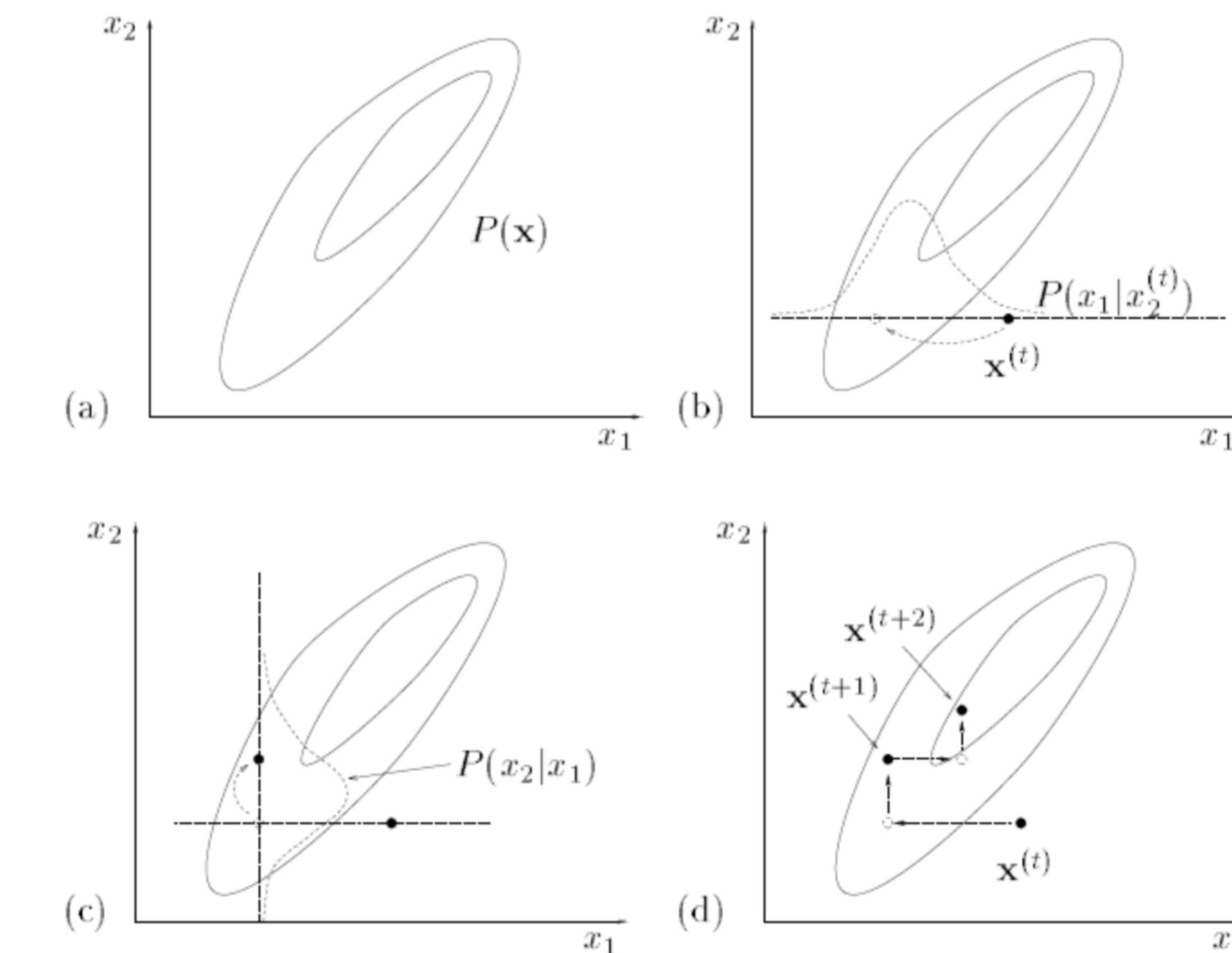
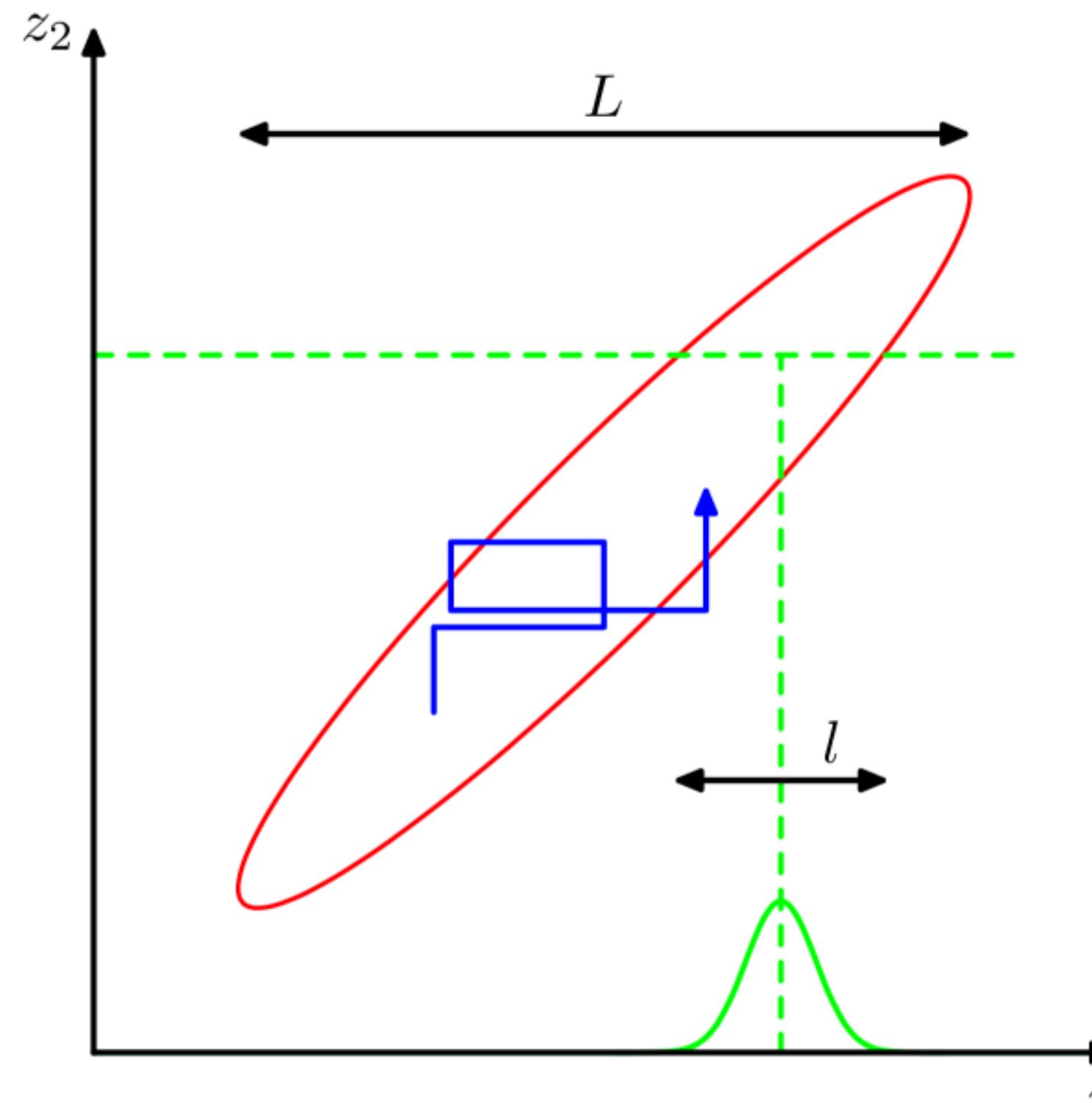


Metropolis chains under high correlation



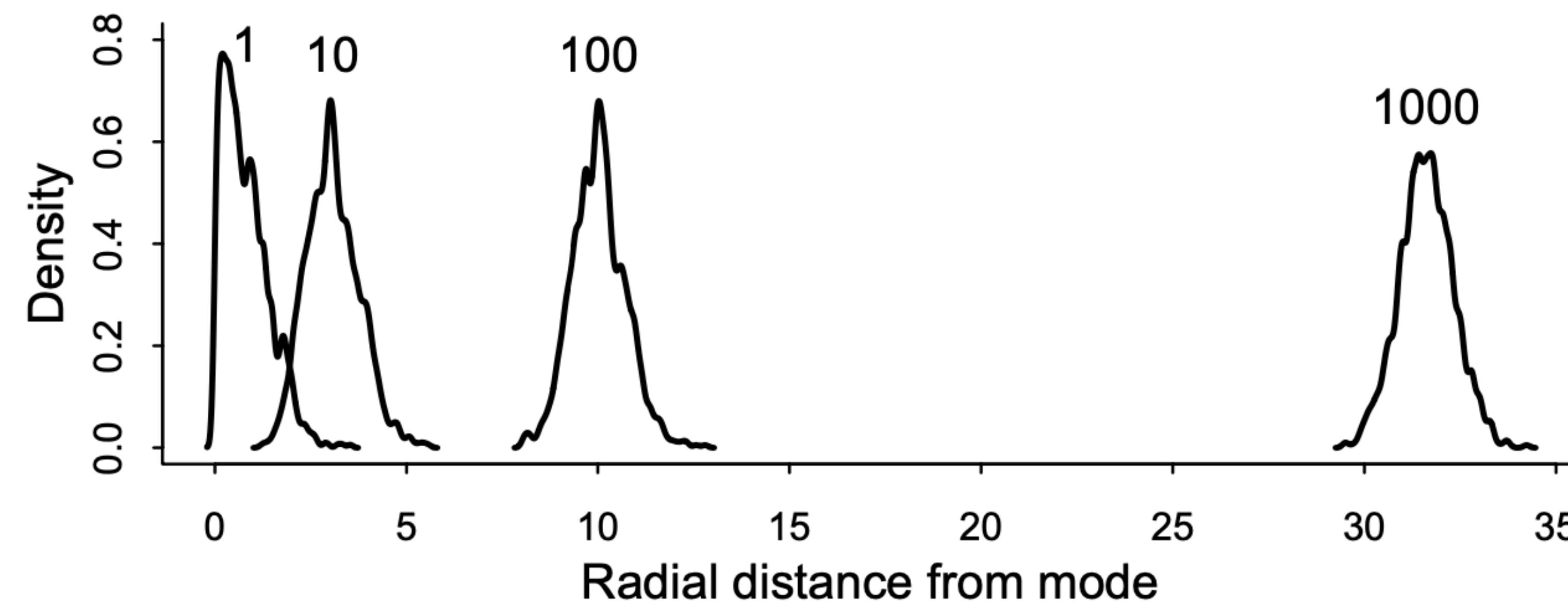
Gibbs Sampling

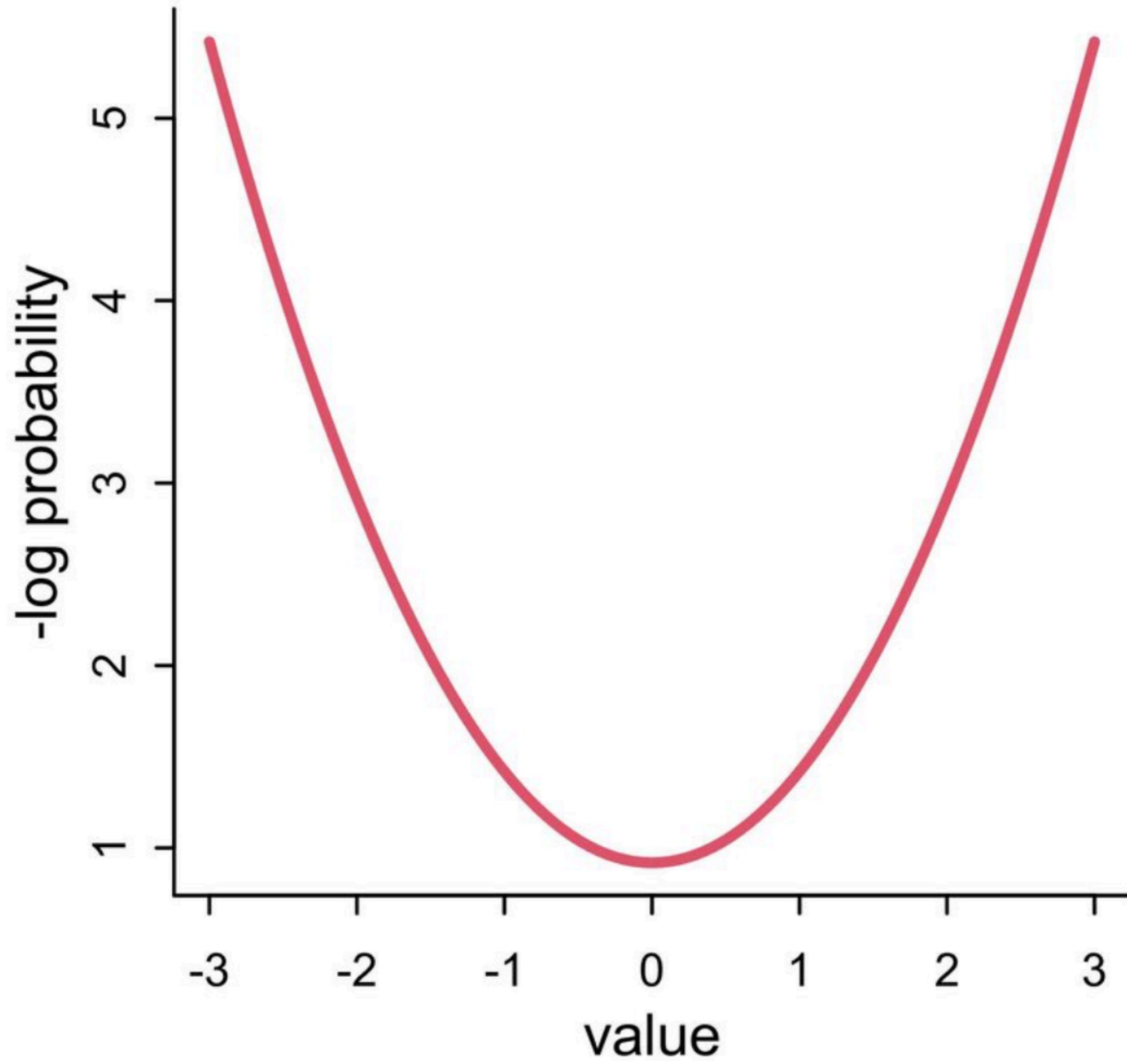
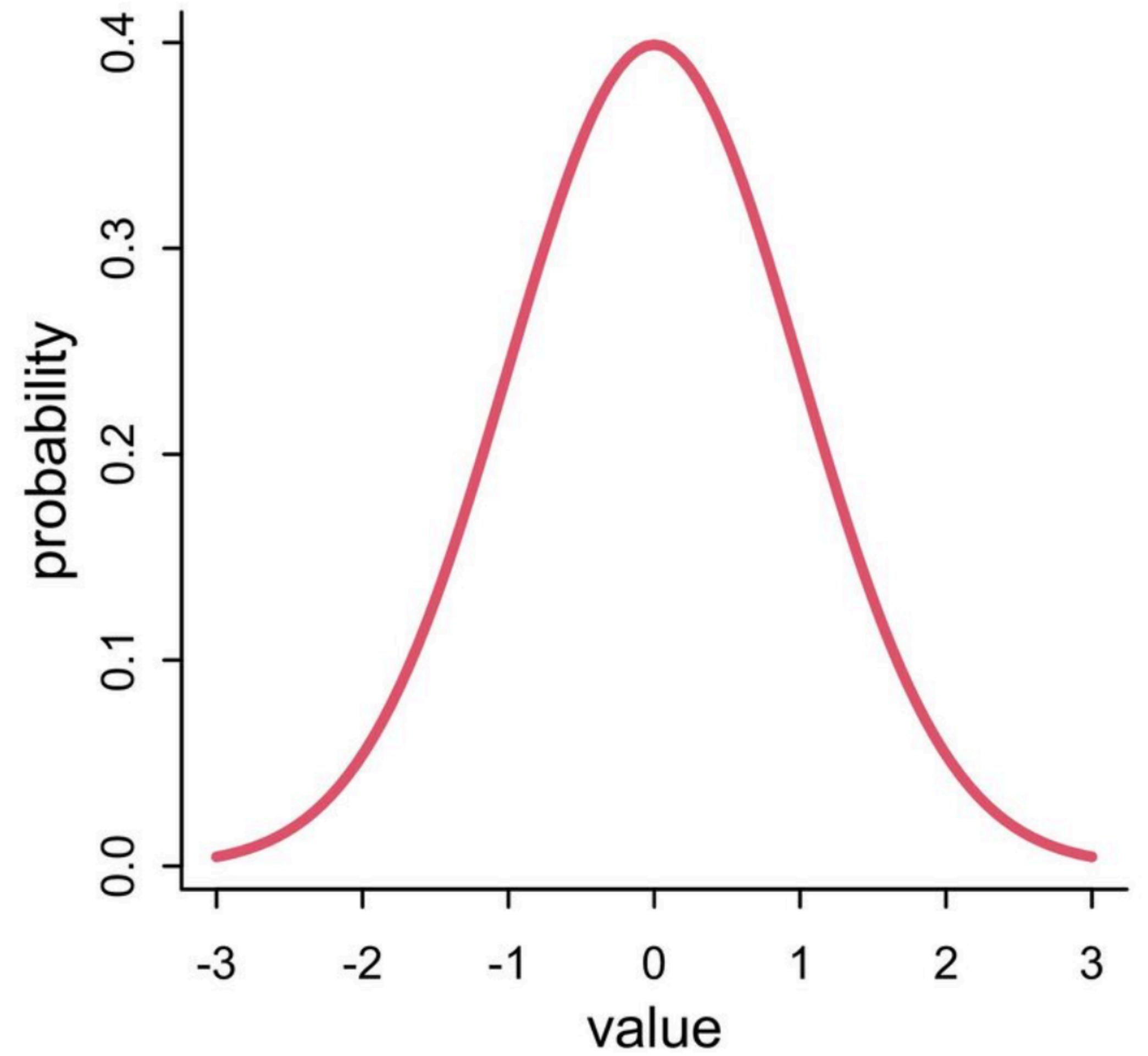
- The idea in Gibbs sampling is that, rather than probabilistically picking the next state all at once, you make a separate probabilistic choice for one of the d dimensions, where each choice depends on the other $d - 1$ dimensions.

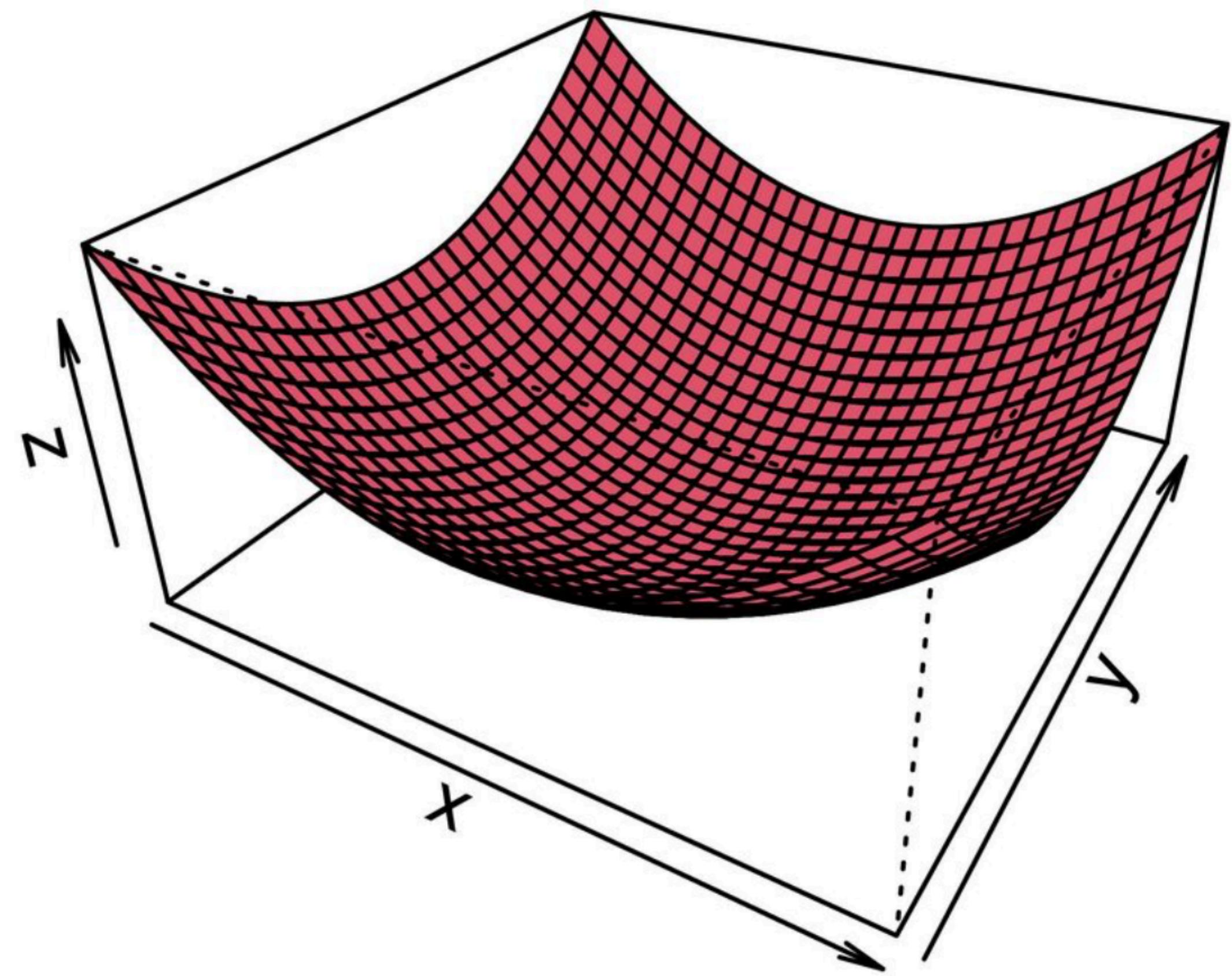
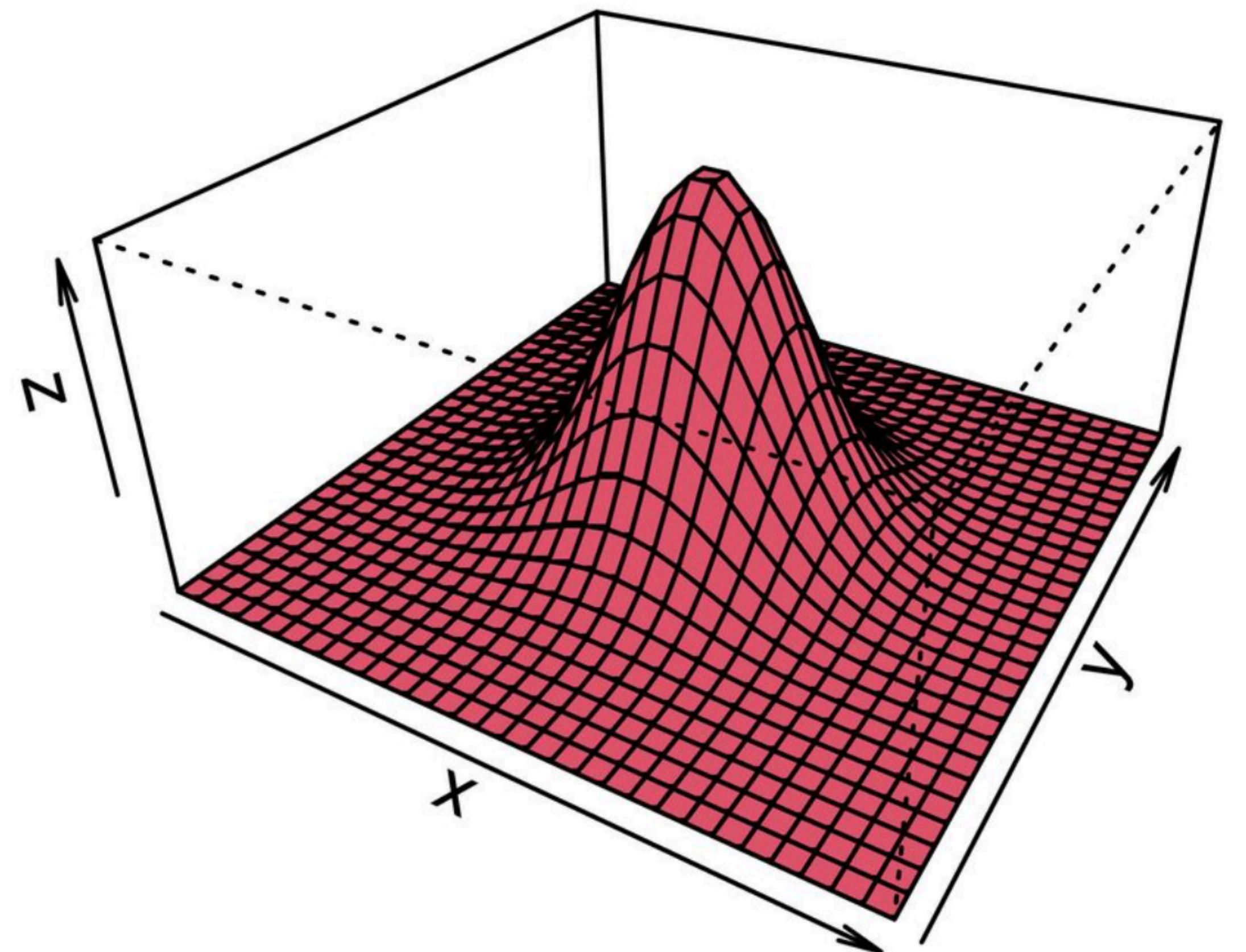


High-dimensional problems. Concentration of measure

- 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





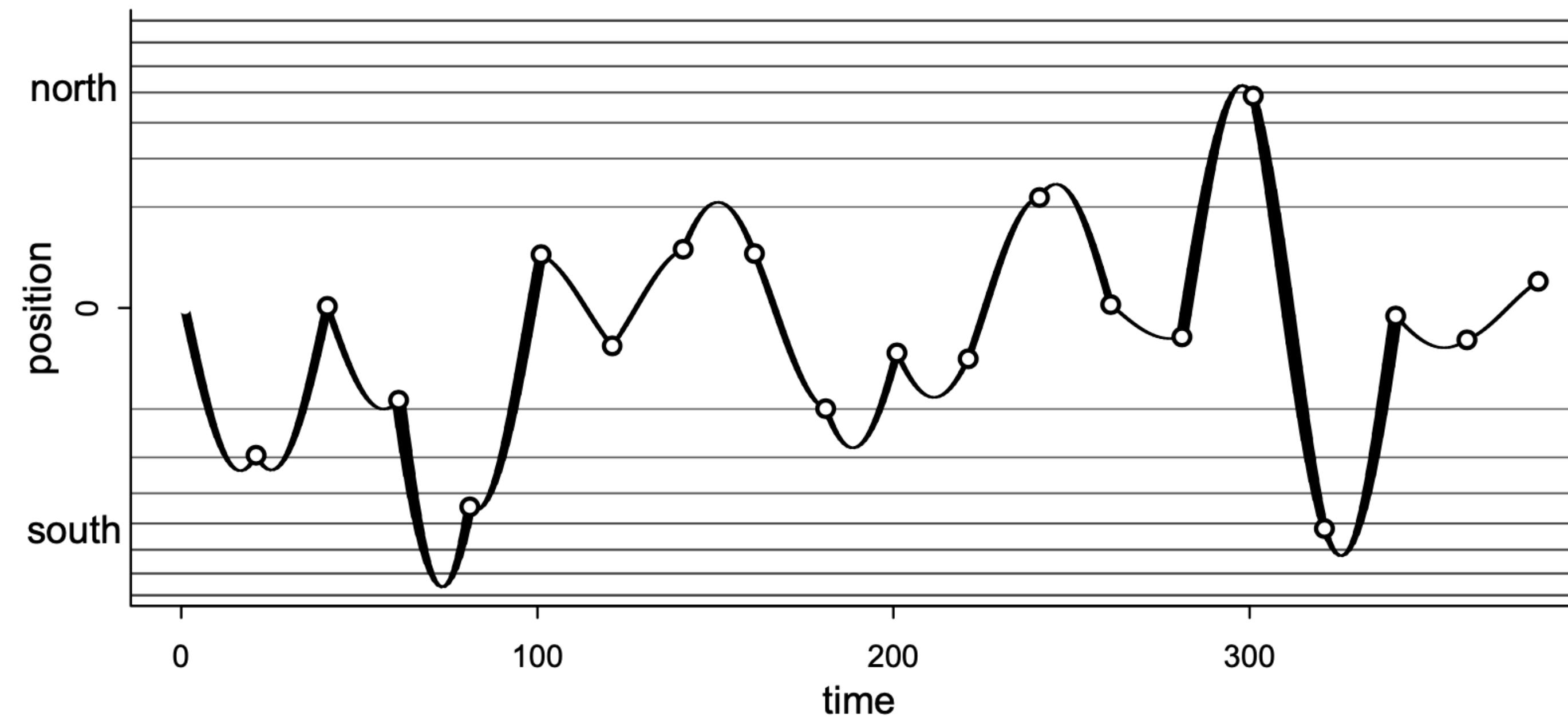






Hamiltonian Monte Carlo (Hybrid Monte Carlo, HMC)

- 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.



Calculus is a superpower

Hamiltonian Monte Carlo needs **gradients**

How does it get them? Write them yourself or...

Auto-diff: Automatic differentiation

Symbolic derivatives of your model code

Used in many machine learning approaches;
“Backpropagation” is special case

$$J = \begin{bmatrix} \frac{\partial B_x}{\partial x} & \frac{\partial B_y}{\partial x} & \frac{\partial B_z}{\partial x} \\ \frac{\partial B_x}{\partial y} & \frac{\partial B_y}{\partial y} & \frac{\partial B_z}{\partial y} \\ \frac{\partial B_x}{\partial z} & \frac{\partial B_y}{\partial z} & \frac{\partial B_z}{\partial z} \end{bmatrix}$$



More calculus)

- Model:

$$x_i \sim \text{Normal}(\mu_x, 1)$$

$$y_i \sim \text{Normal}(\mu_y, 1)$$

$$\mu_x \sim \text{Normal}(0, 0.5)$$

$$\mu_y \sim \text{Normal}(0, 0.5)$$

- Where the following formula came from?

$$\sum_i \log p(y_i | \mu_y, 1) + \sum_i \log p(x_i | \mu_x, 1) + \log p(\mu_y | 0, 0.5) + \log p(\mu_x, 0, 0.5)$$

Gradients, 2D Gaussian example

Log posterior:

$$\sum_i \log p(y_i|\mu_y, 1) + \sum_i \log p(x_i|\mu_x, 1) + \log p(\mu_y|0, 0.5) + \log p(\mu_x|0, 0.5)$$

- Given we have a Gaussian density, $N(y|a,b)$, then

$$\frac{\partial \log N(y|a,b)}{\partial a} = \frac{y-a}{b^2}$$

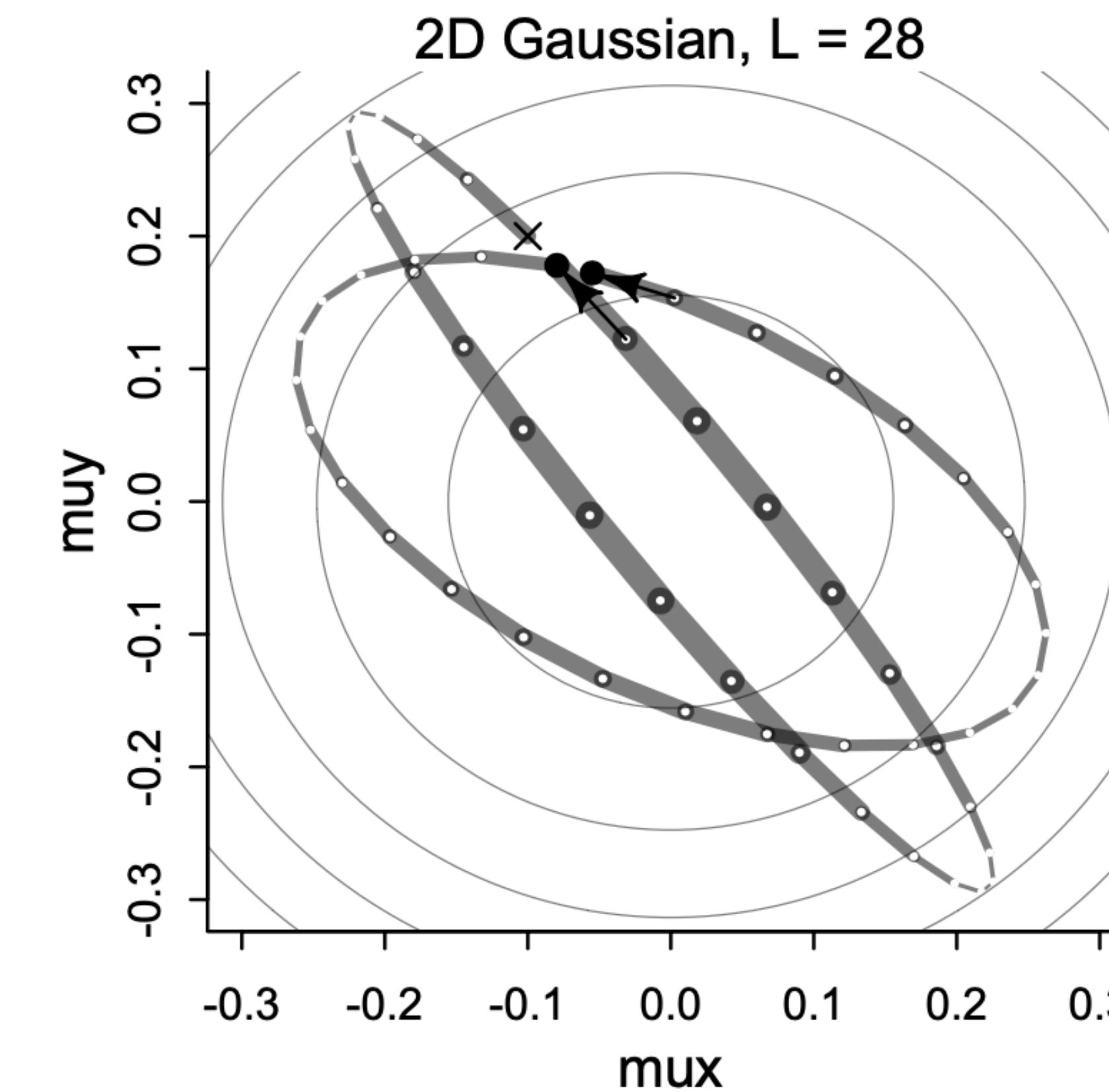
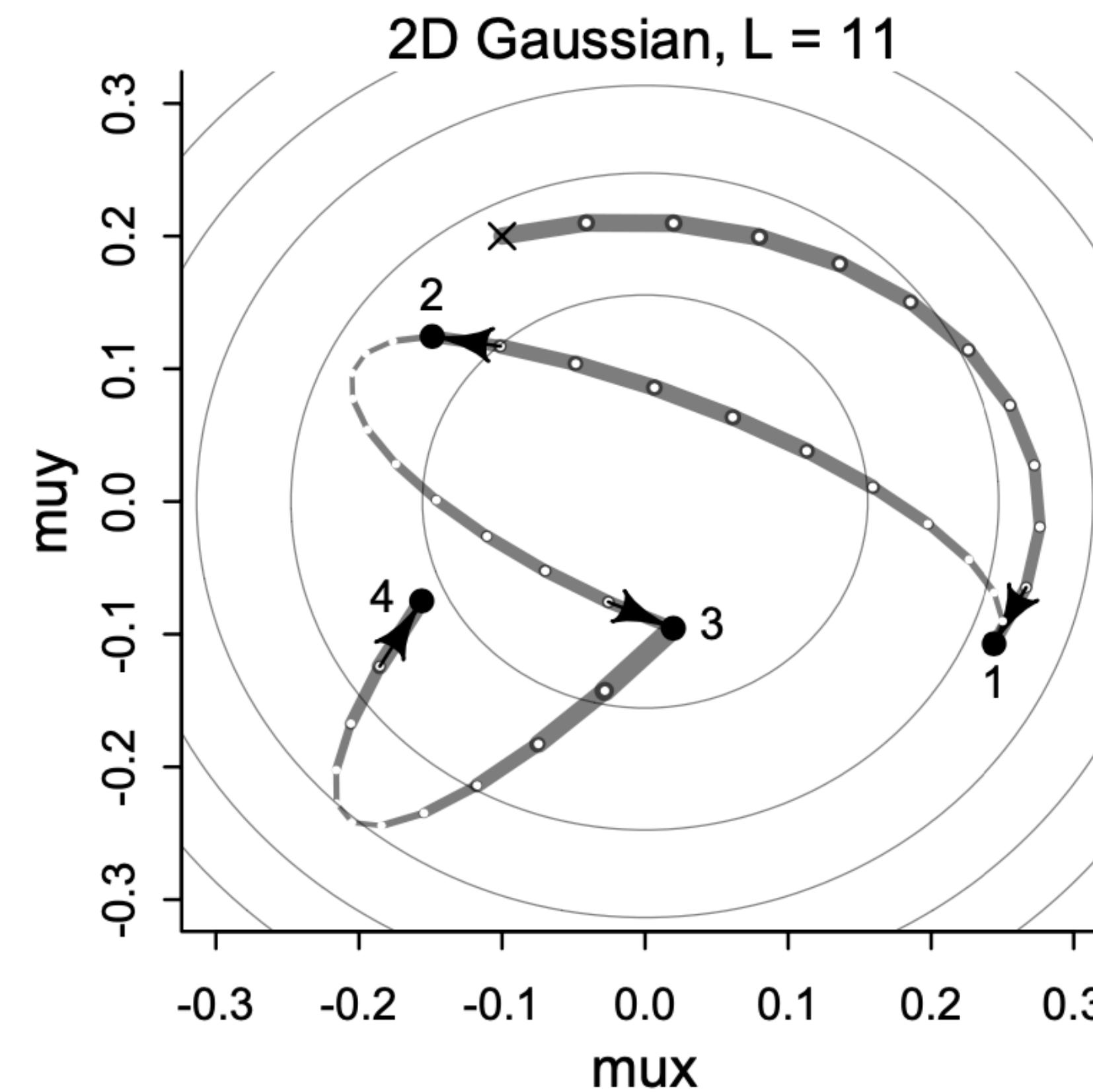
- for first dim.:

$$\frac{\partial U}{\partial \mu_x} = \frac{\partial \log N(x|\mu_x, 1)}{\partial \mu_x} + \frac{\partial \log N(\mu_x|0, 0.5)}{\partial \mu_x} = \sum_i \frac{x_i - \mu_x}{1^2} + \frac{0 - \mu_x}{0.5^2}$$

R Code

```
# gradient function
# need vector of partial derivatives of U with respect to vector q
U_gradient <- function( q , a=0 , b=1 , k=0 , d=1 ) {
  muy <- q[1]
  mux <- q[2]
  G1 <- sum( y - muy ) + (a - muy)/b^2 #dU/dmuy
  G2 <- sum( x - mux ) + (k - mux)/d^2 #dU/dmux
  return( c( -G1 , -G2 ) ) # negative bc energy is neg-log-prob
}
# test data
set.seed(7)
y <- rnorm(50)
x <- rnorm(50)
x <- as.numeric(scale(x))
y <- as.numeric(scale(y))
```

U-turns



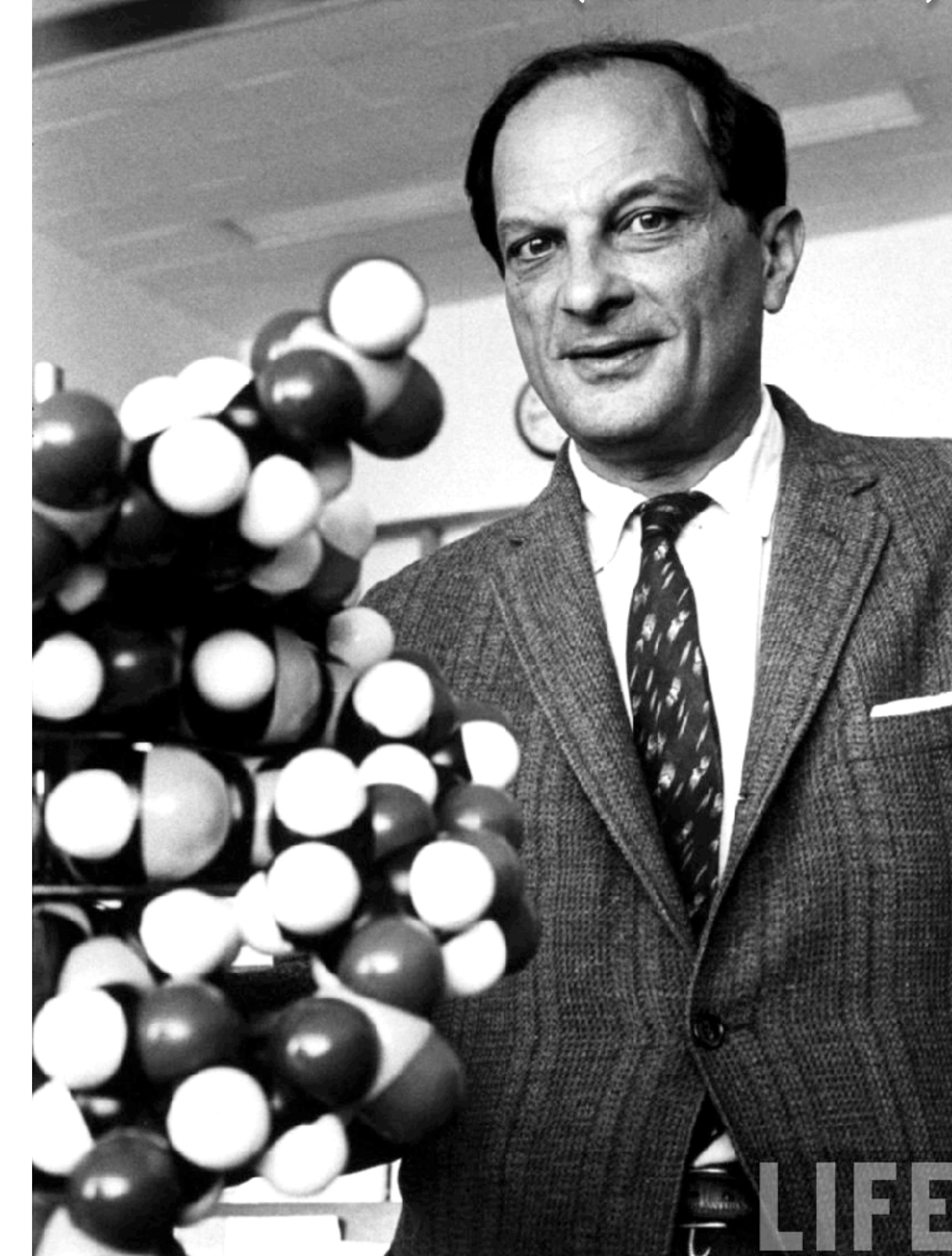


Stan

About Stan

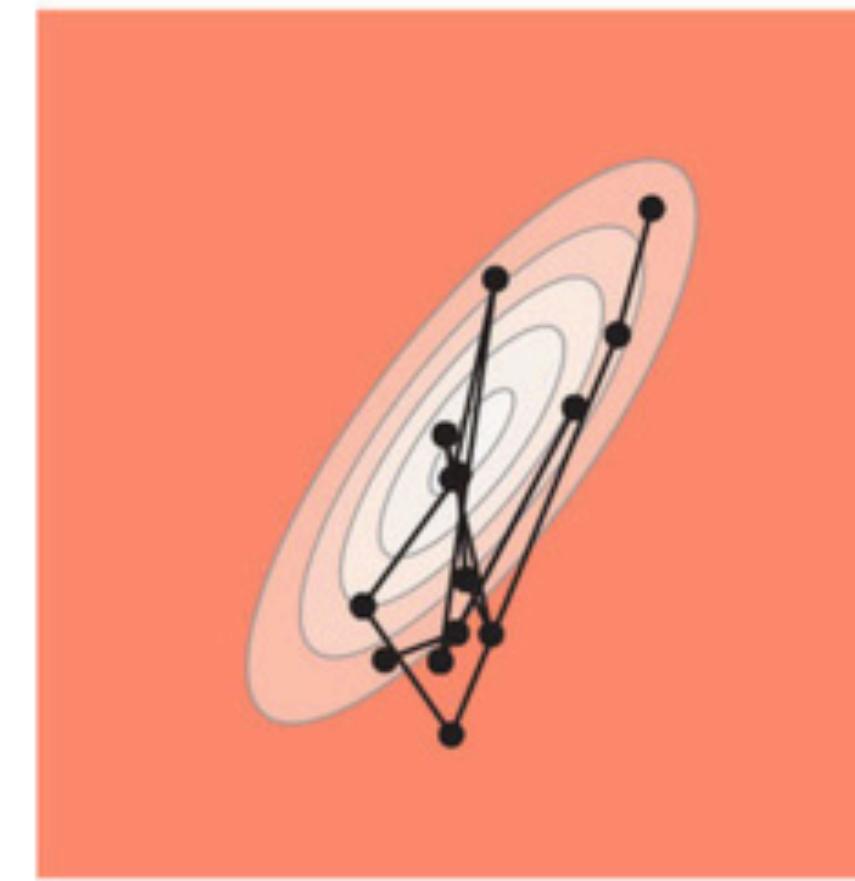
Stan is a state-of-the-art platform for statistical modeling and high-performance statistical computation. Thousands of users rely on Stan for statistical modeling, data analysis, and prediction in the social, biological, and physical sciences, engineering, and business.

Stanisław Ulam (1909–1984)

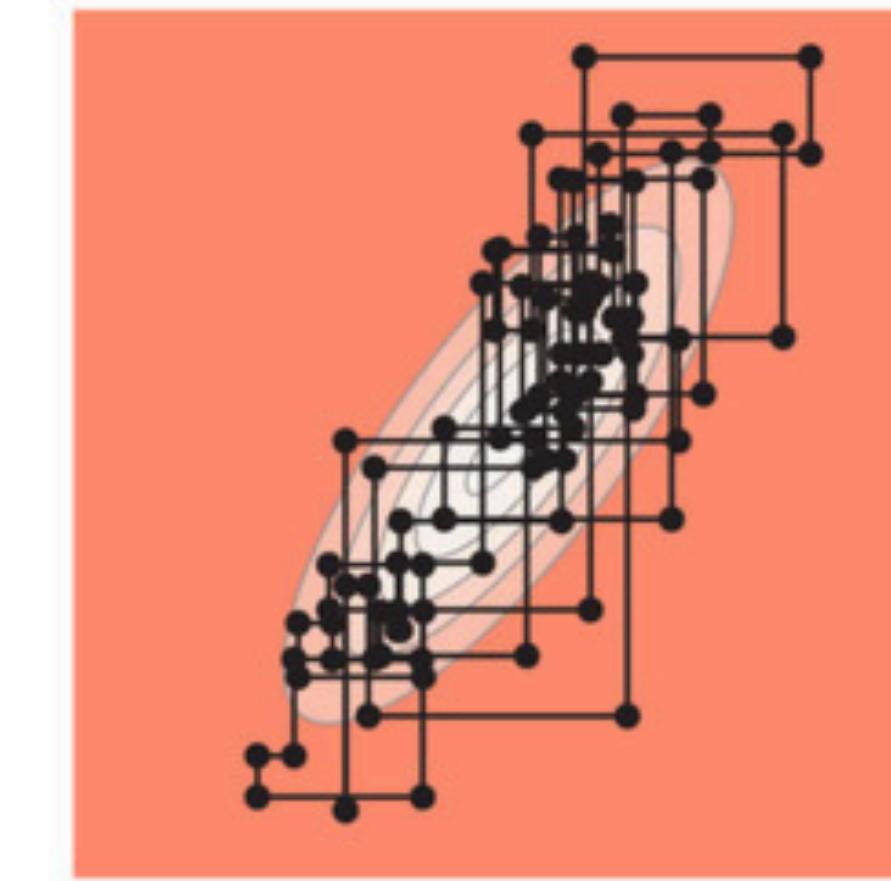


Different Sampling methods

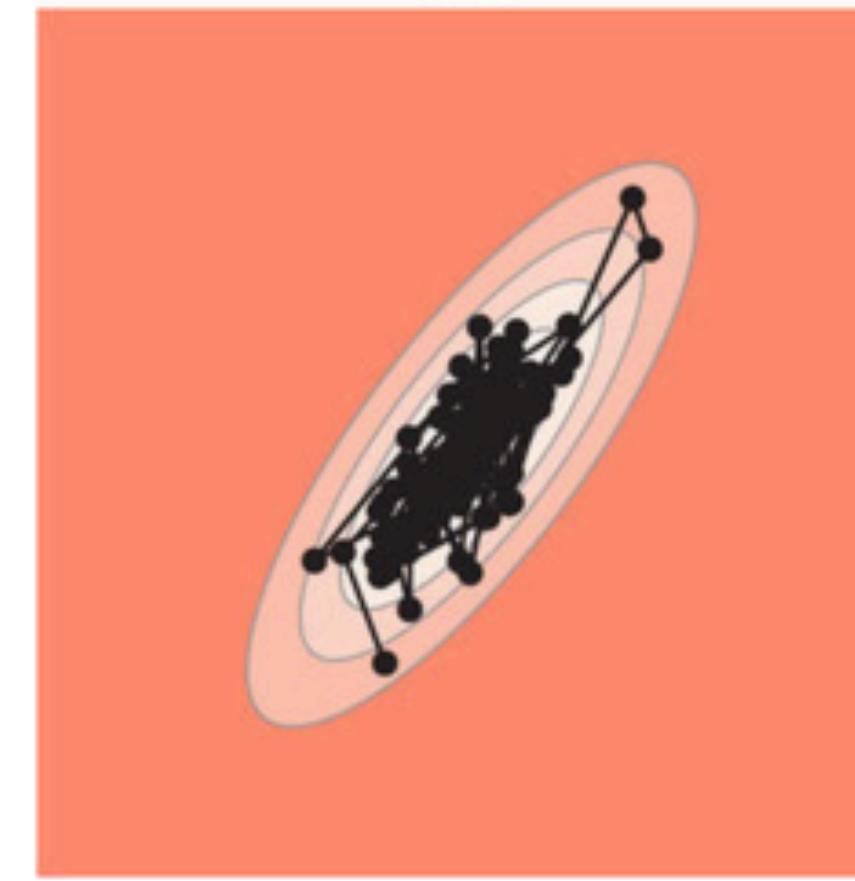
Random Walk Metropolis



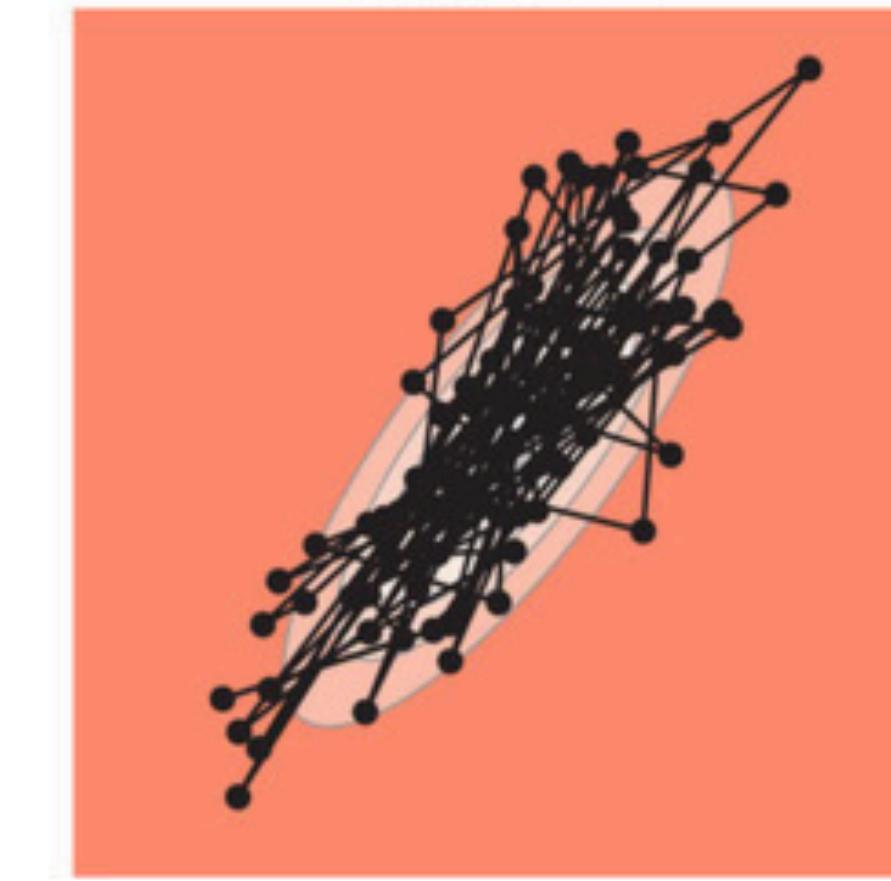
Gibbs



HMC

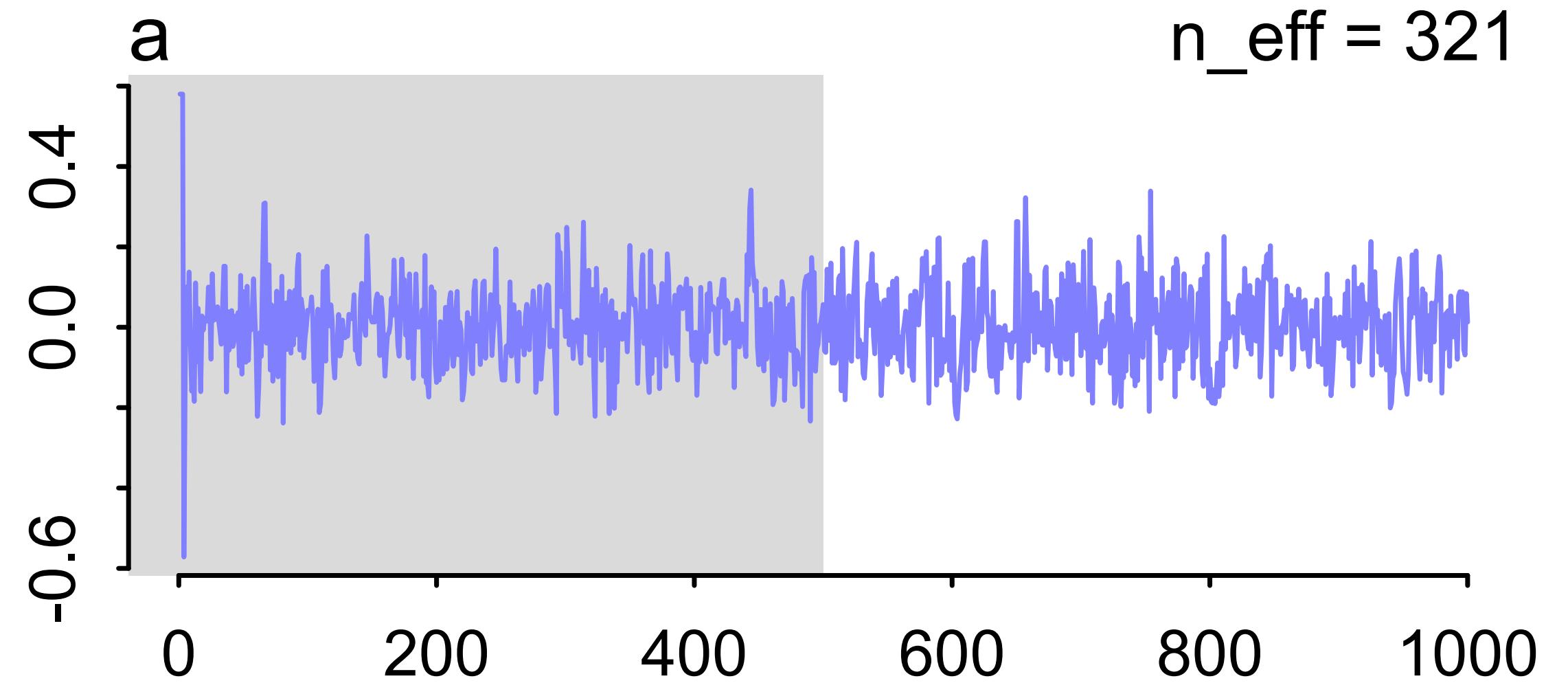


Independent

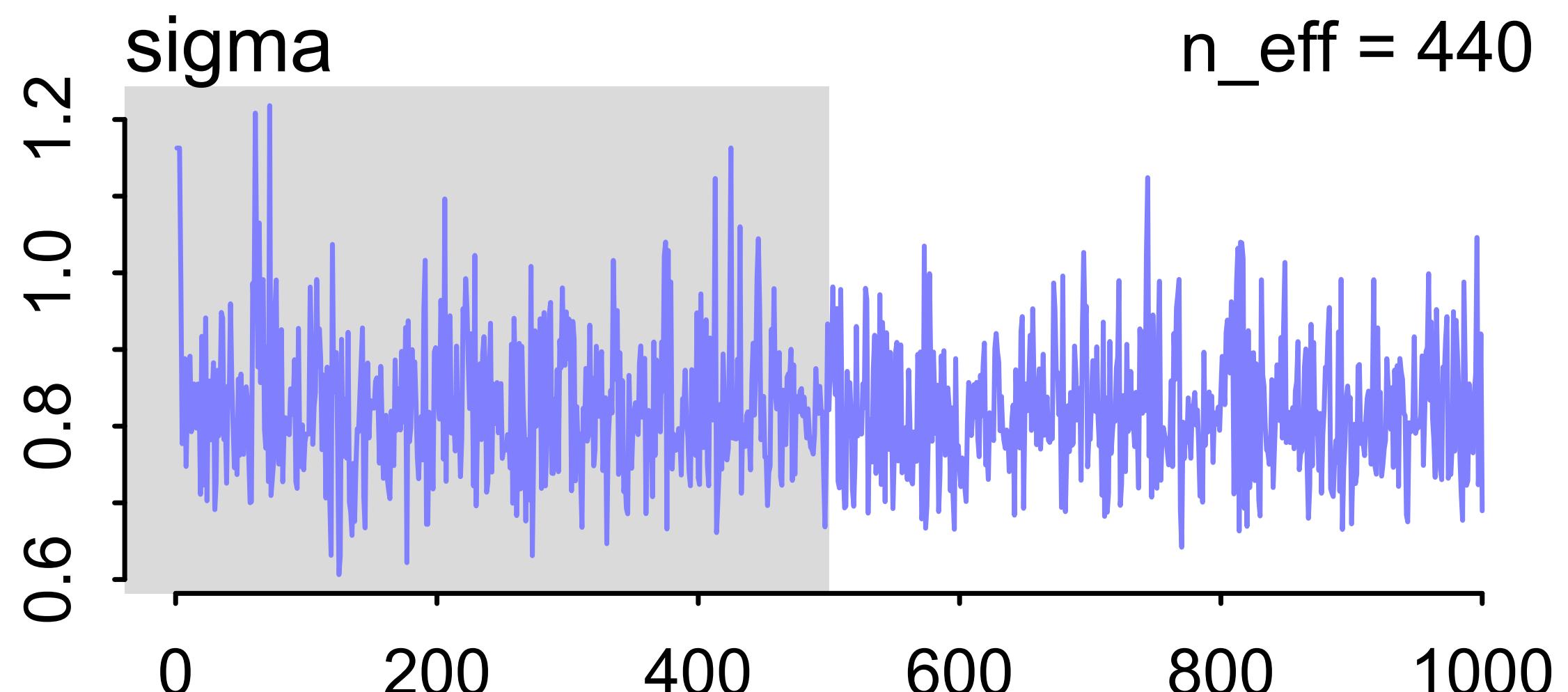
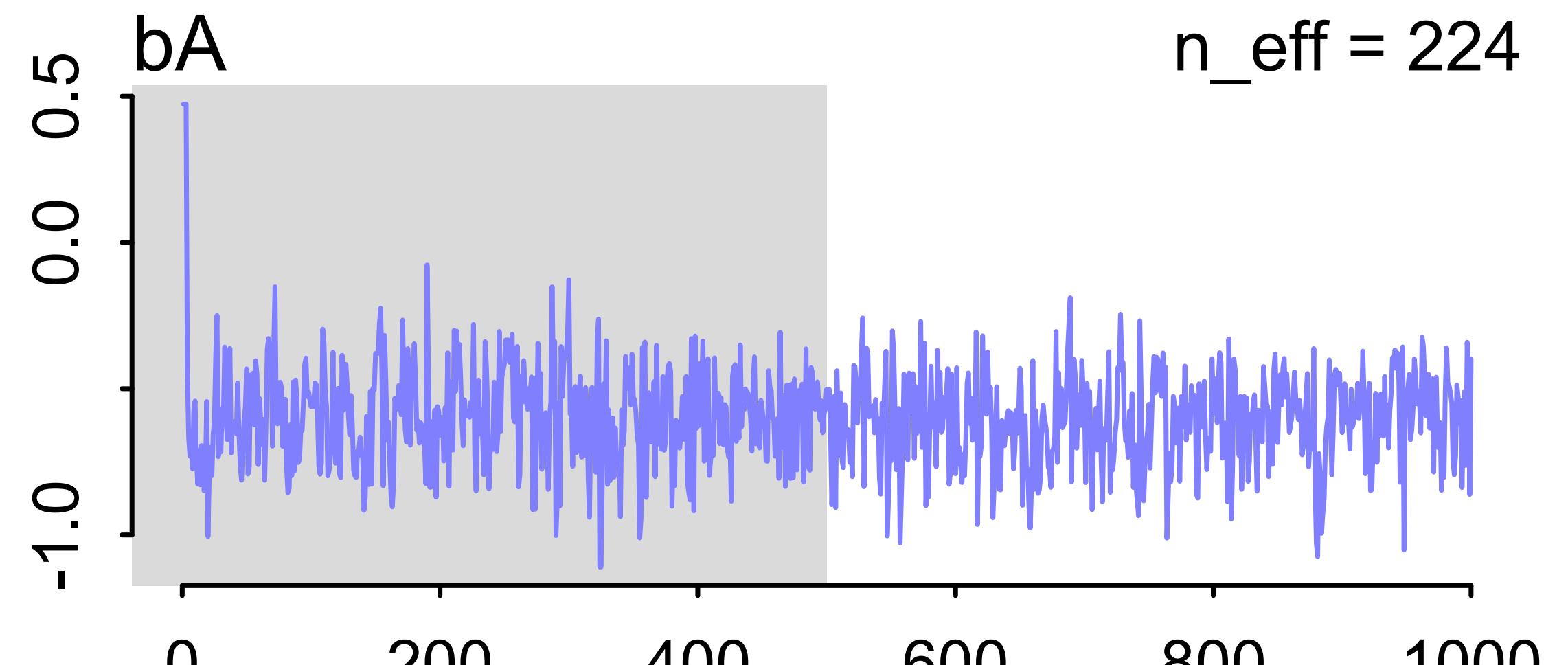
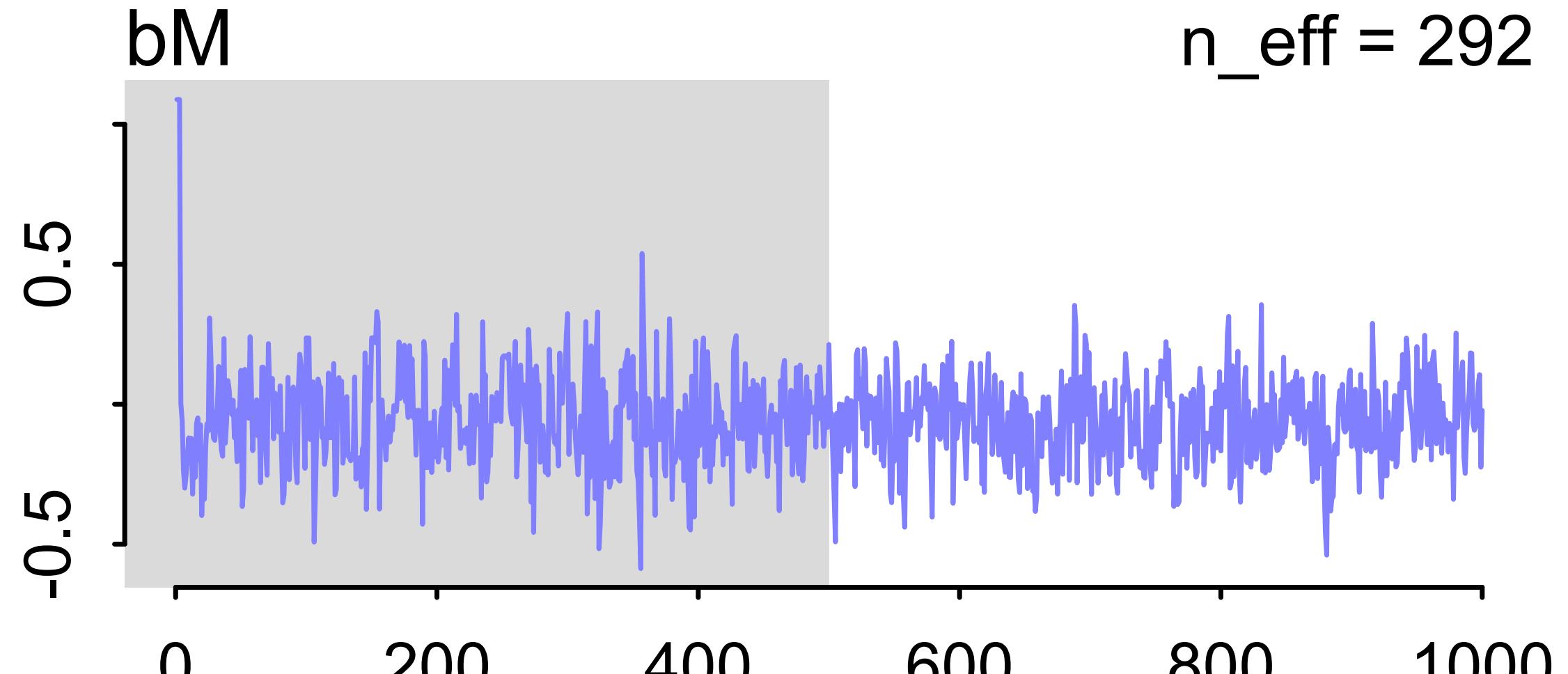
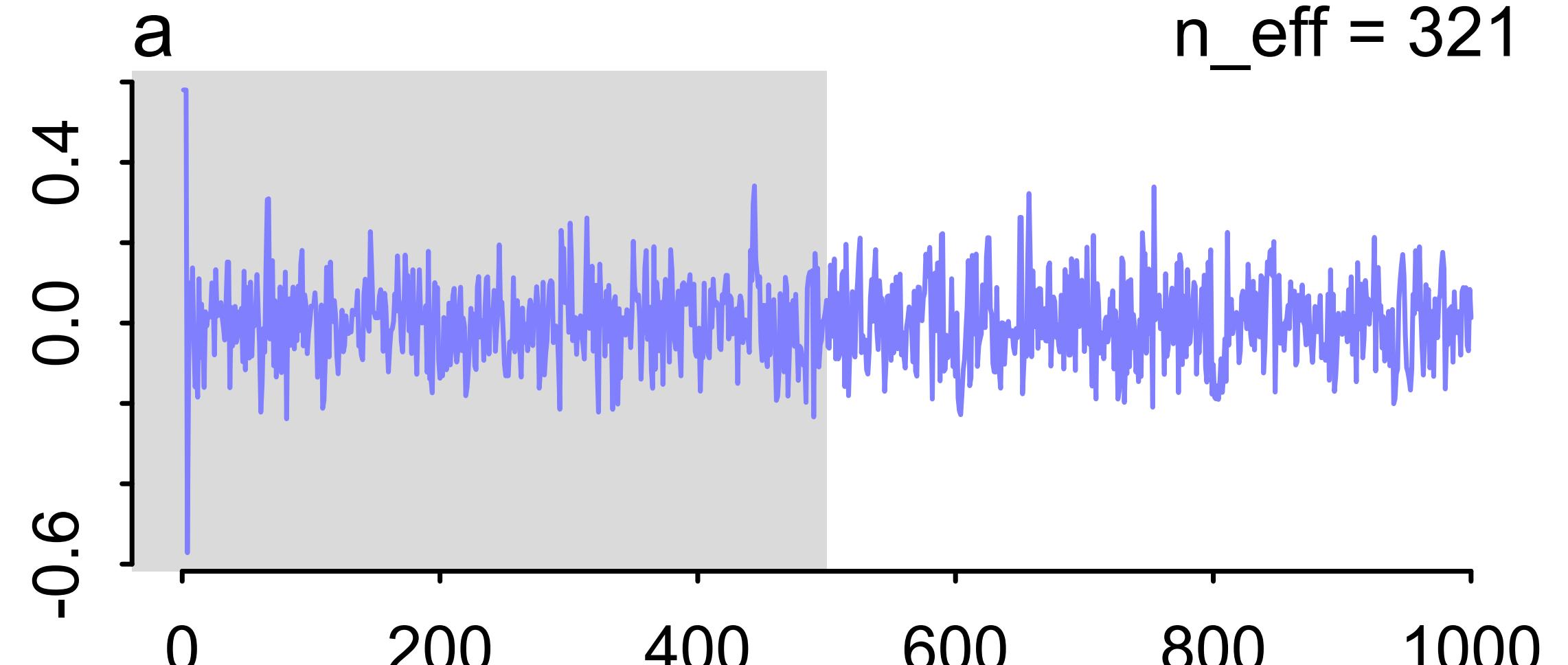


- Lambert, Ben. A Students Guide to Bayesian Statistics SAGE Publications. 2018

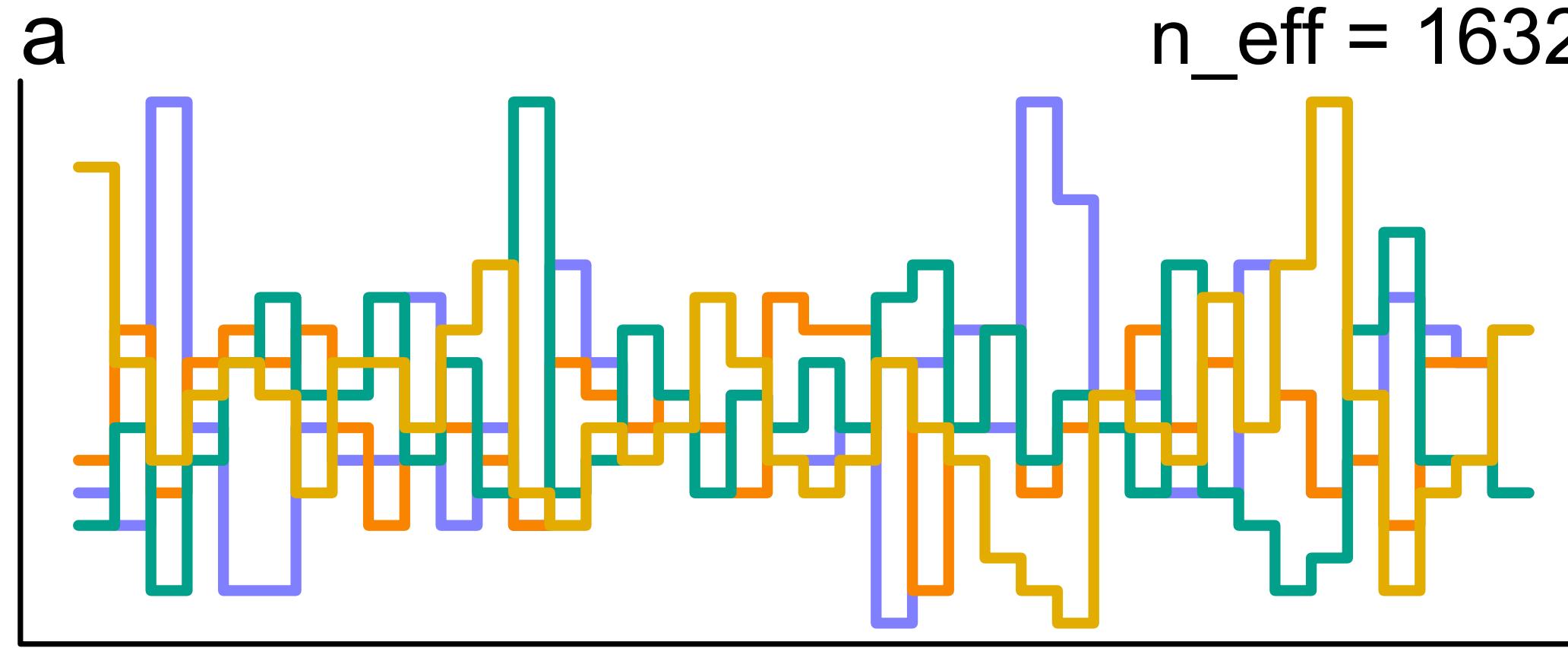
Trace plots



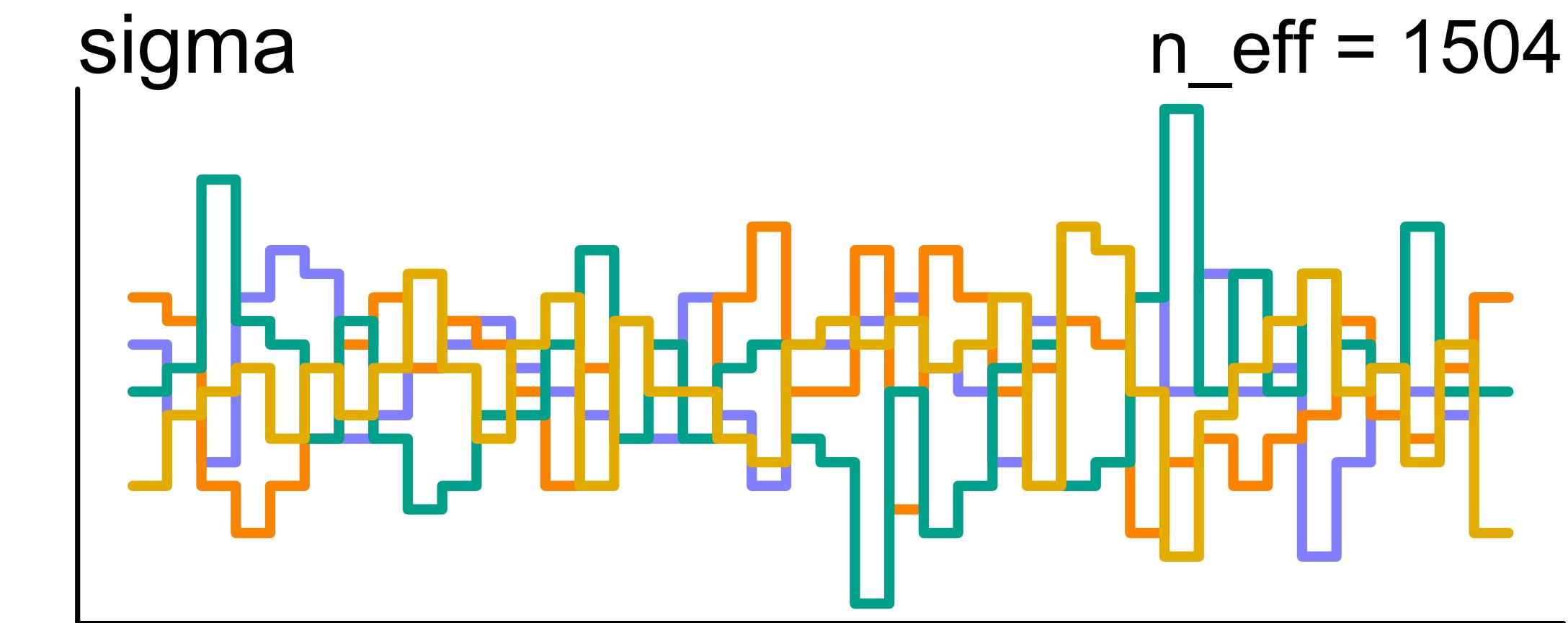
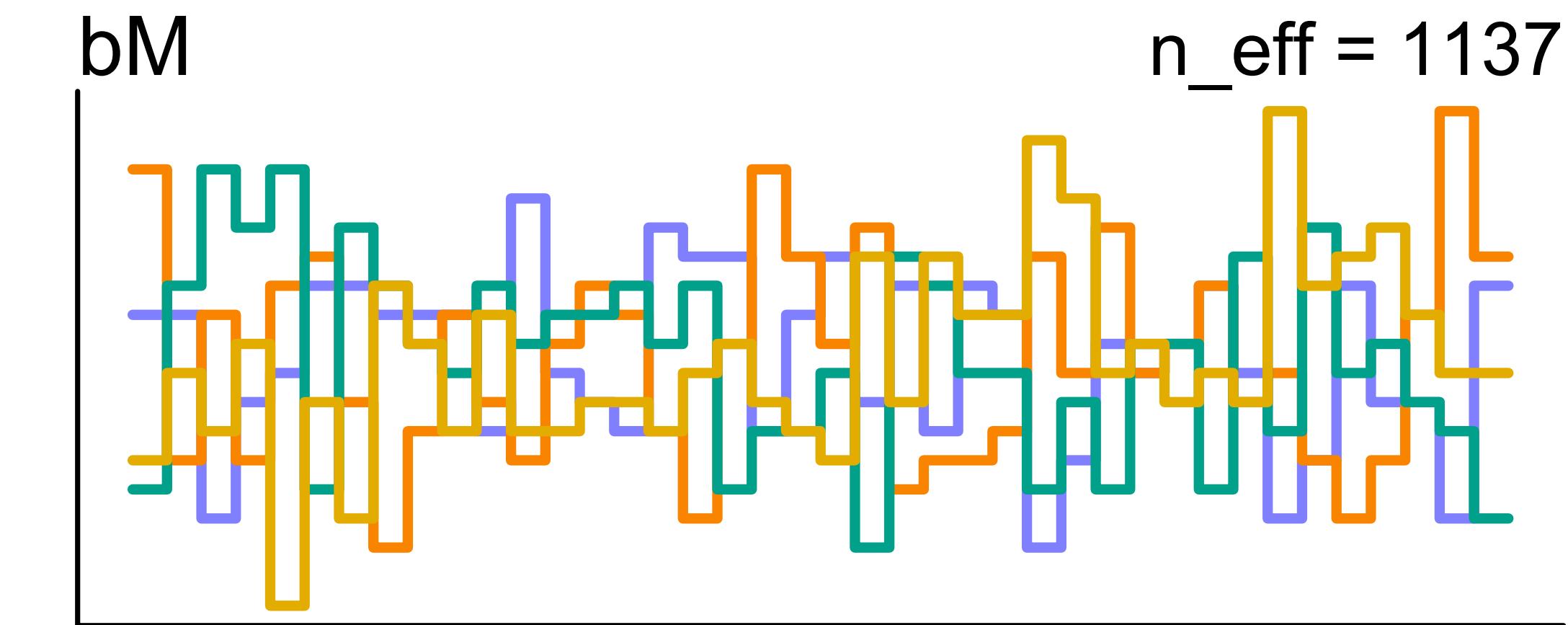
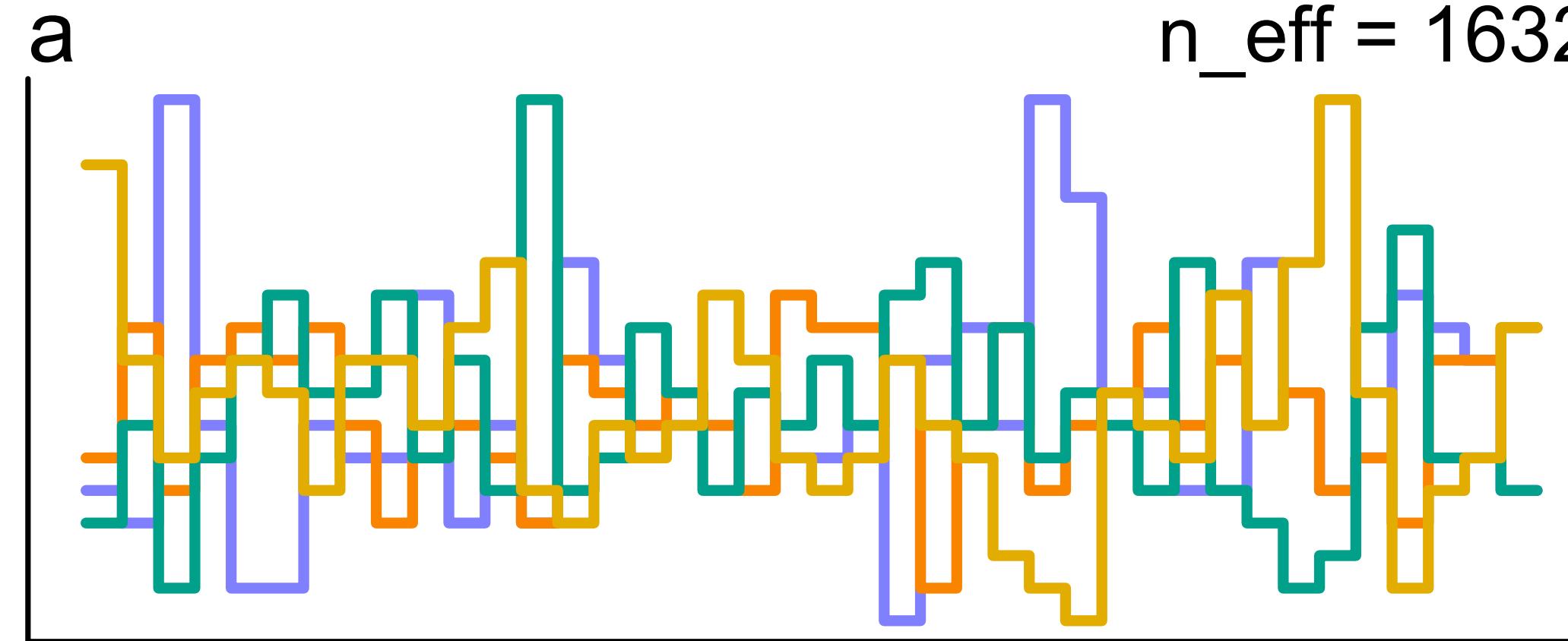
Trace plots

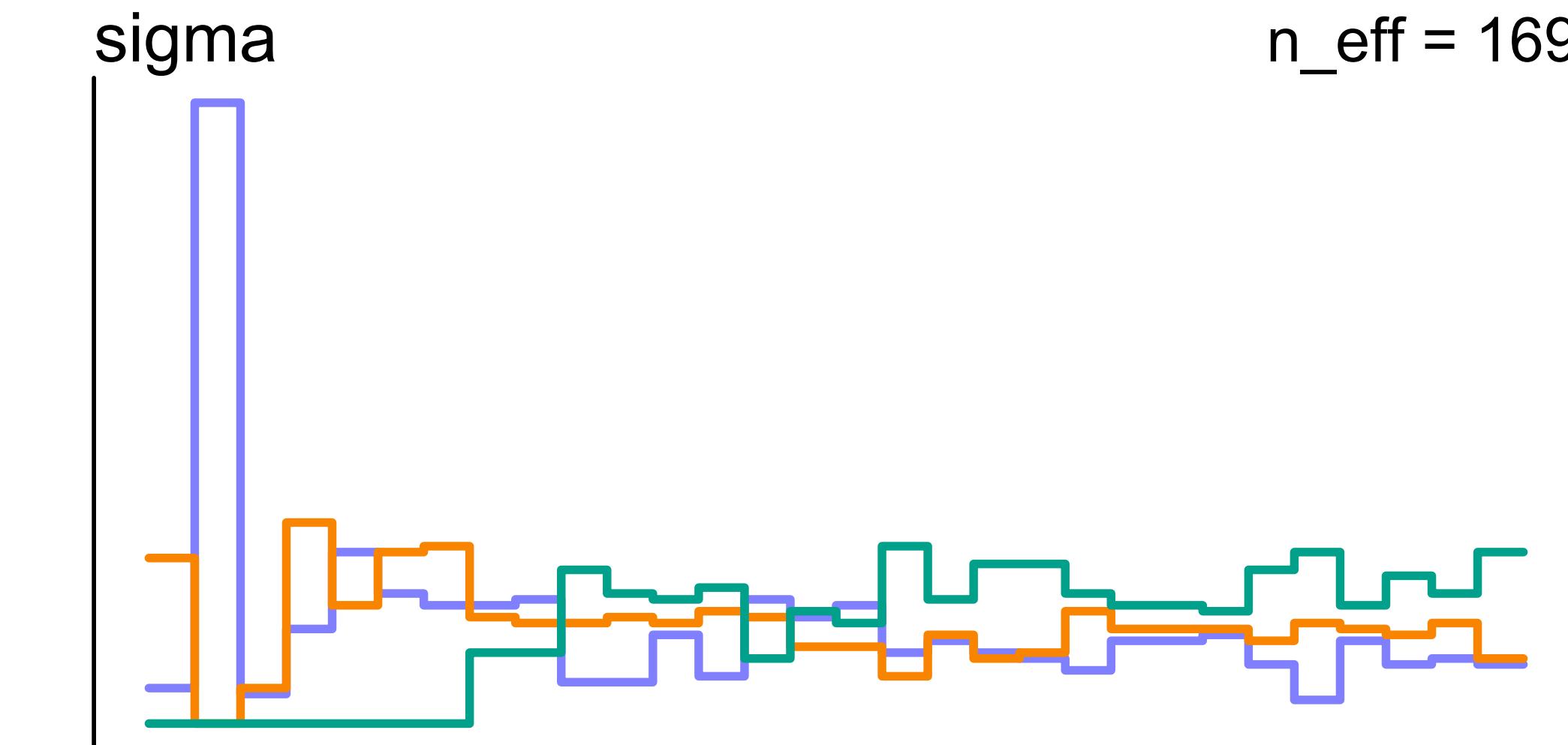
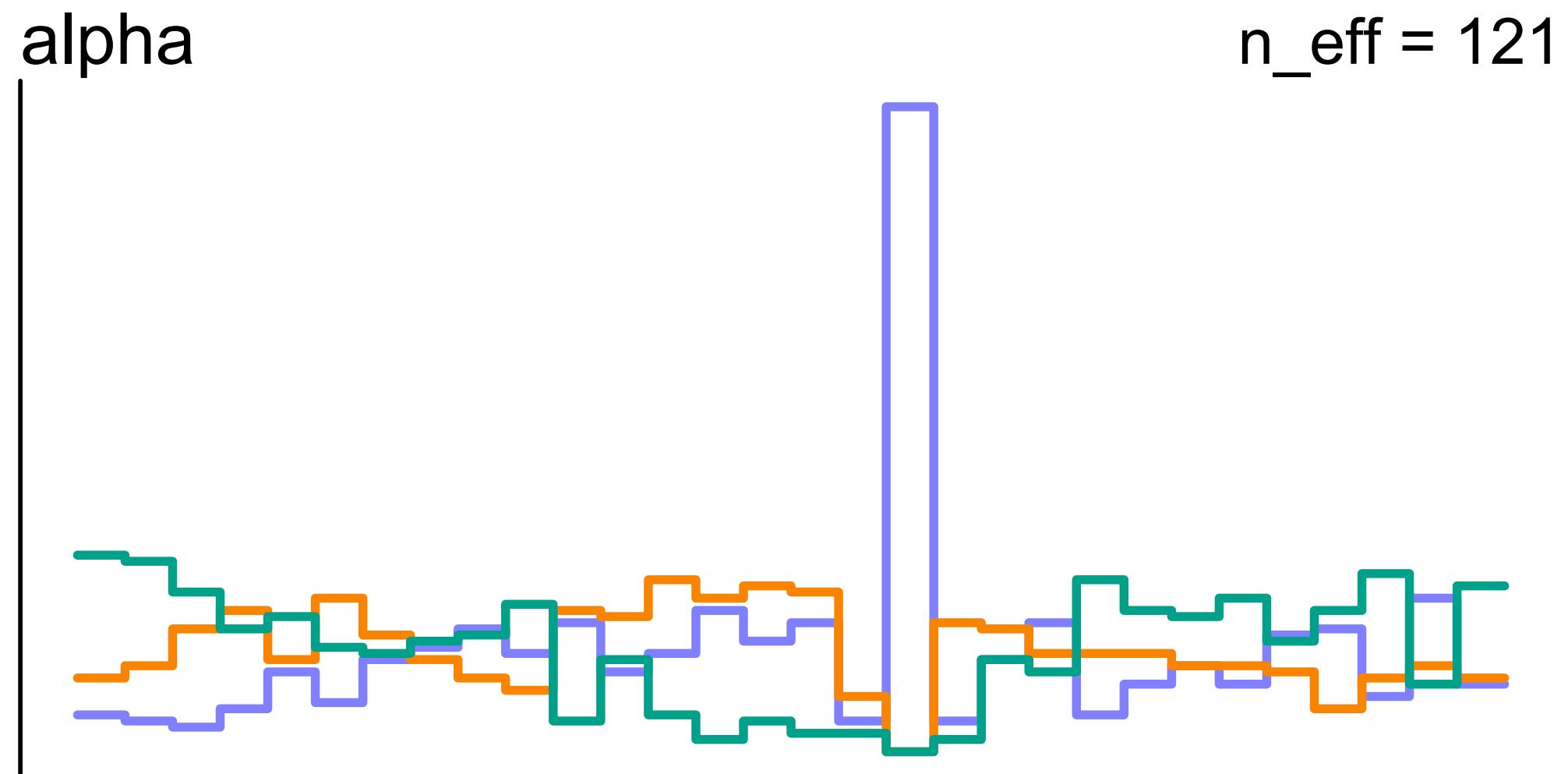
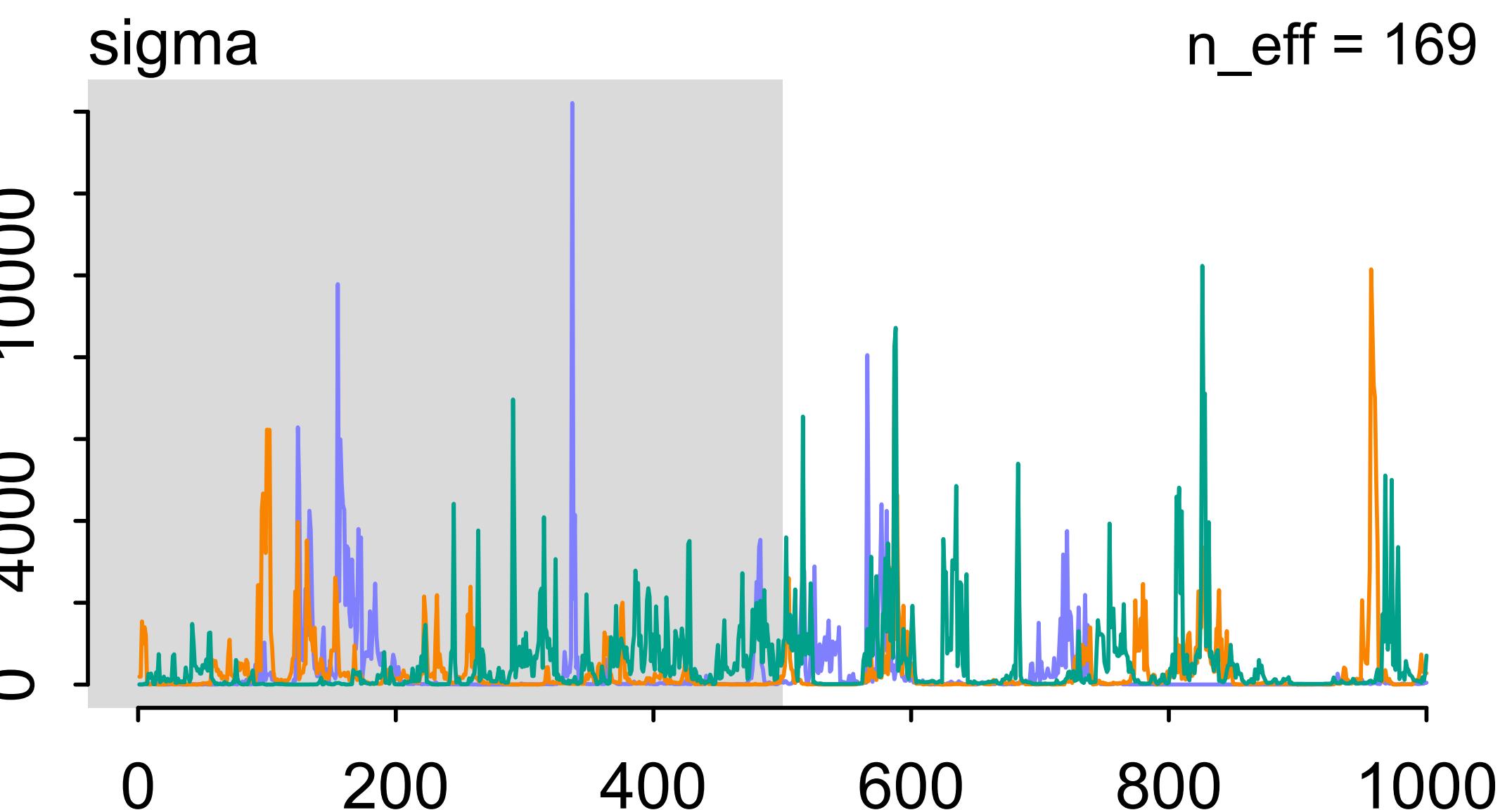
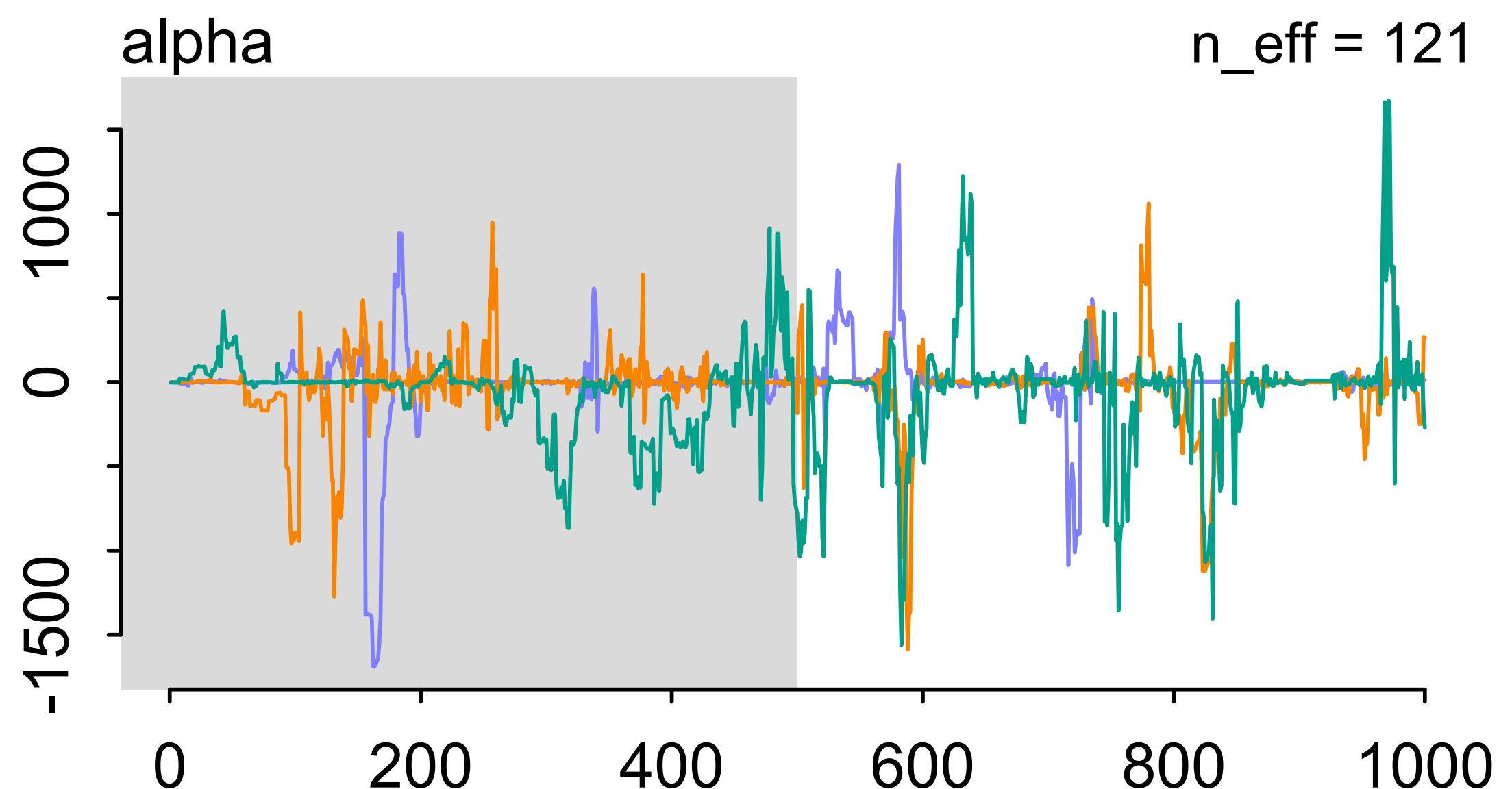


Trace rank (Trank) plots



Trace rank (Trank) plots





R-hat

When chains converge:

(1) Start and end of each chain explores same region

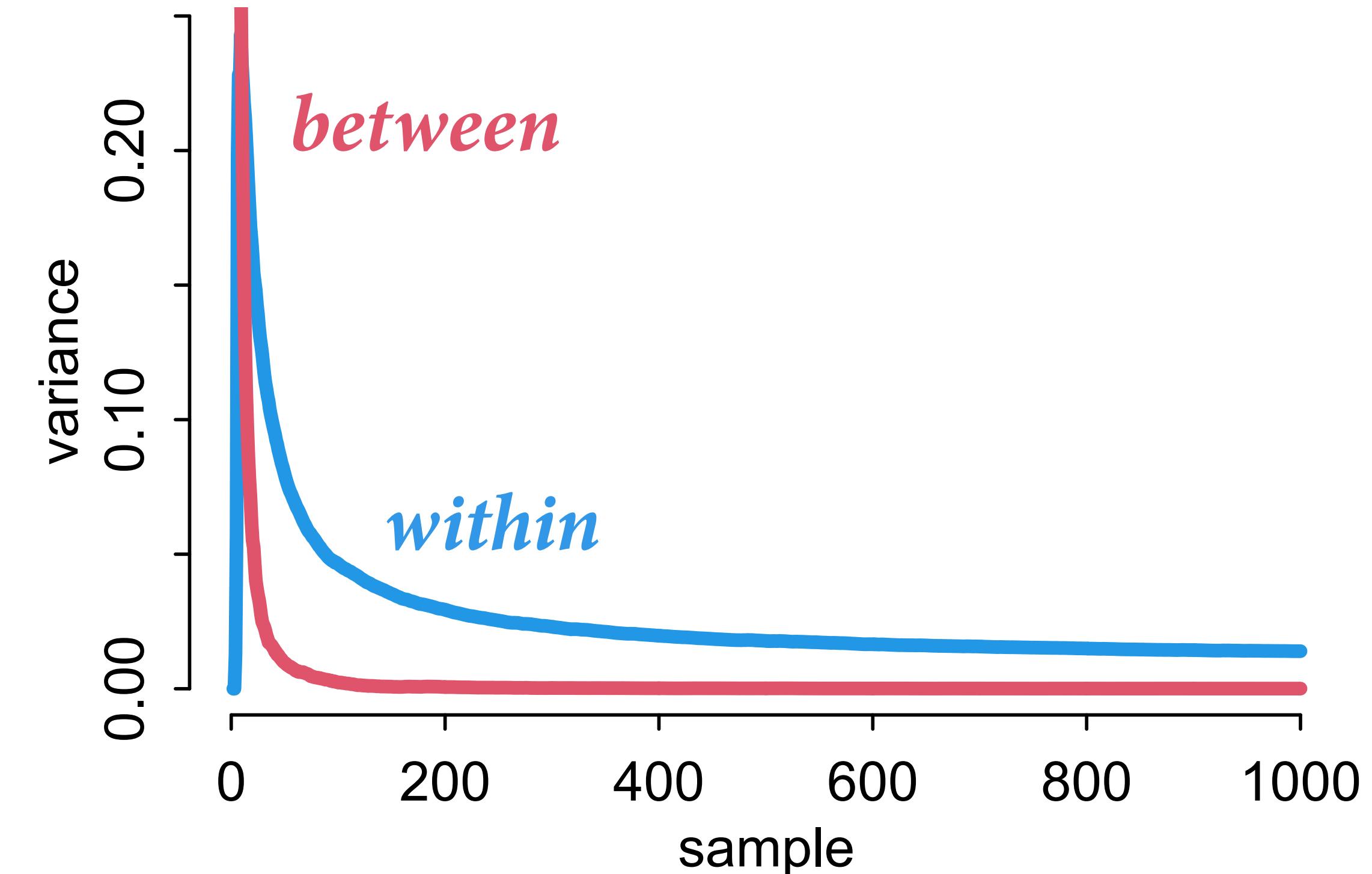
(2) Independent chains explore same region

R-hat is a ratio of variances:

As total variance shrinks to average variance within chains, R-hat approaches 1

NO GUARANTEES; NOT A TEST

```
> precis(mHMC)
    mean   sd  5.5% 94.5% n_eff Rhat4
a     0.00 0.10 -0.16  0.16  1632    1
bM    -0.06 0.17 -0.32  0.21  1137    1
bA    -0.61 0.17 -0.86 -0.34  1160    1
sigma 0.83 0.09  0.70  0.99  1504    1
>
```



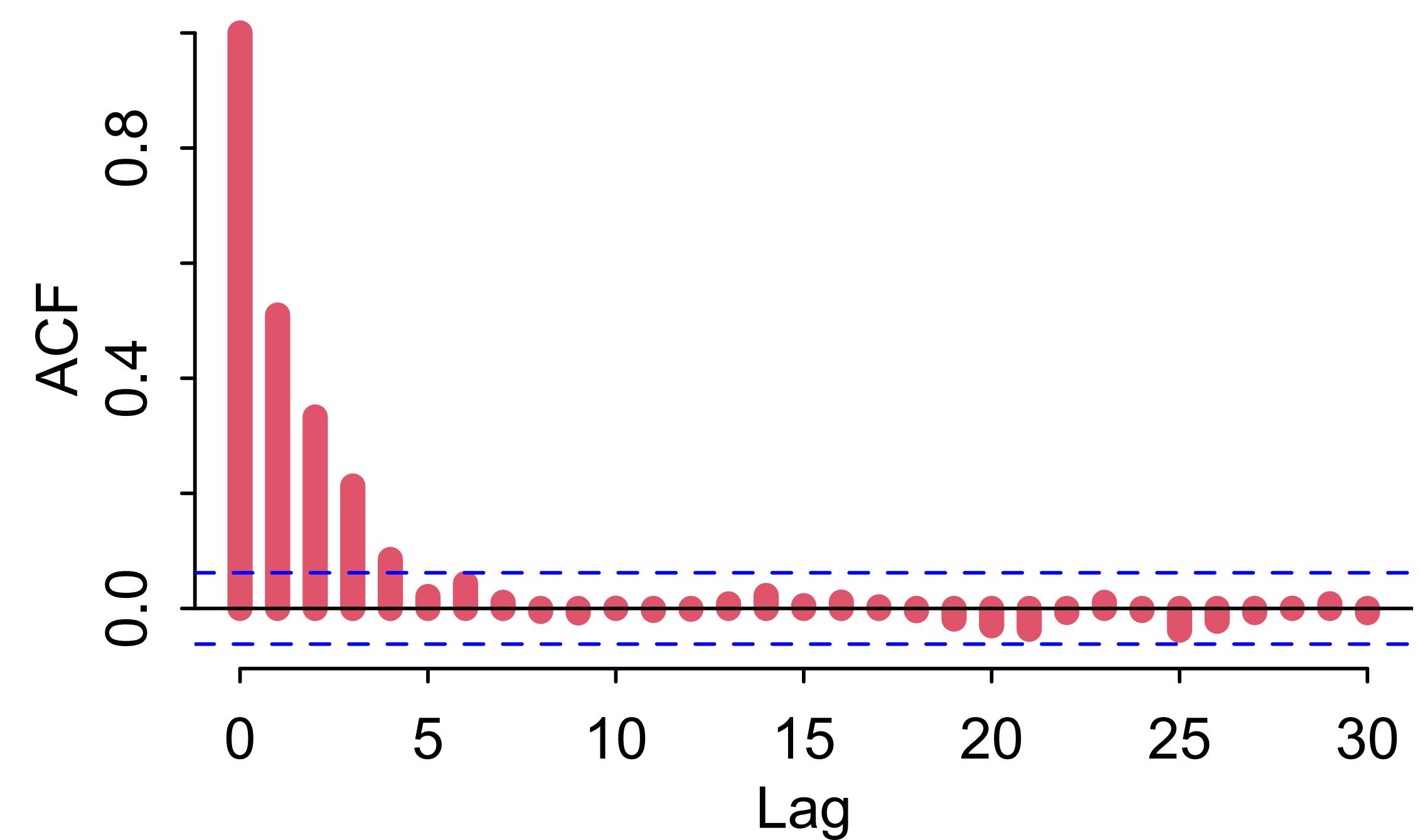
n_eff

Estimate of number of **effective samples**

“How long would the chain be, if each sample was independent of the one before it?”

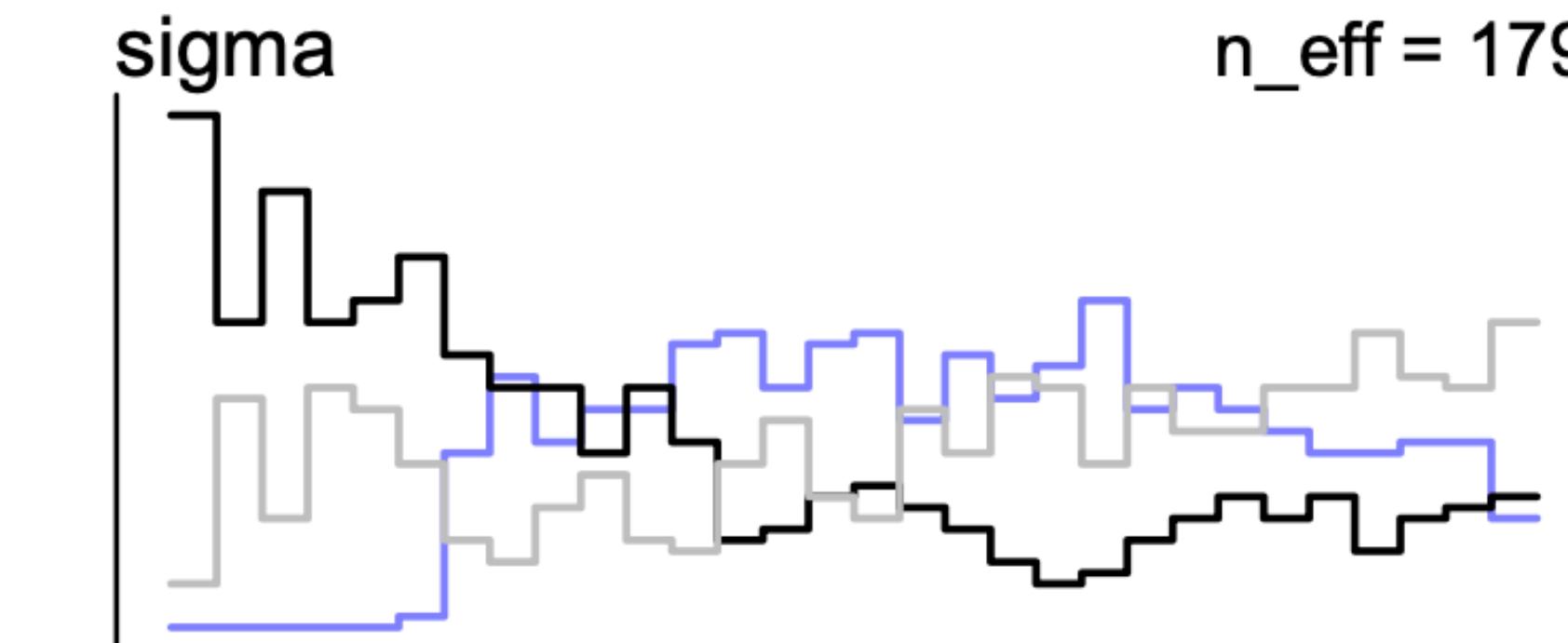
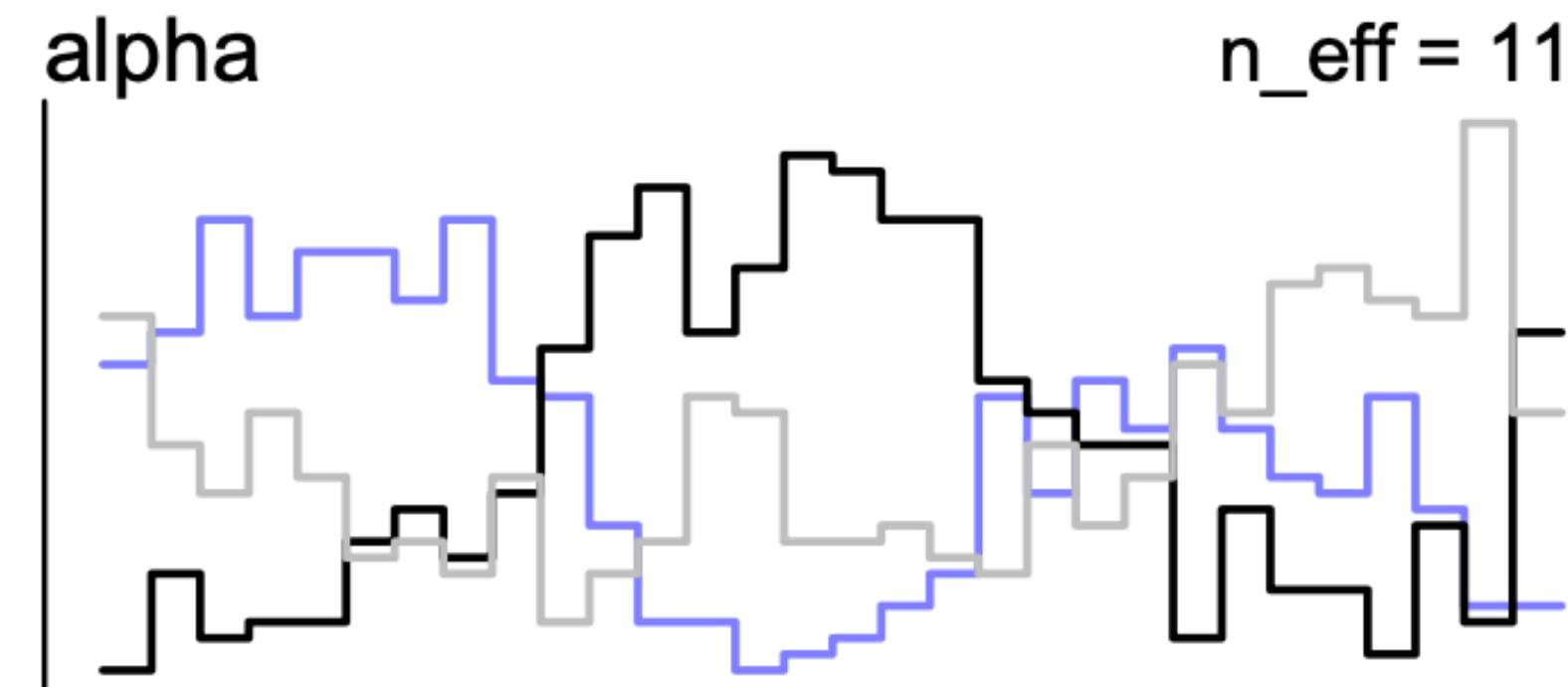
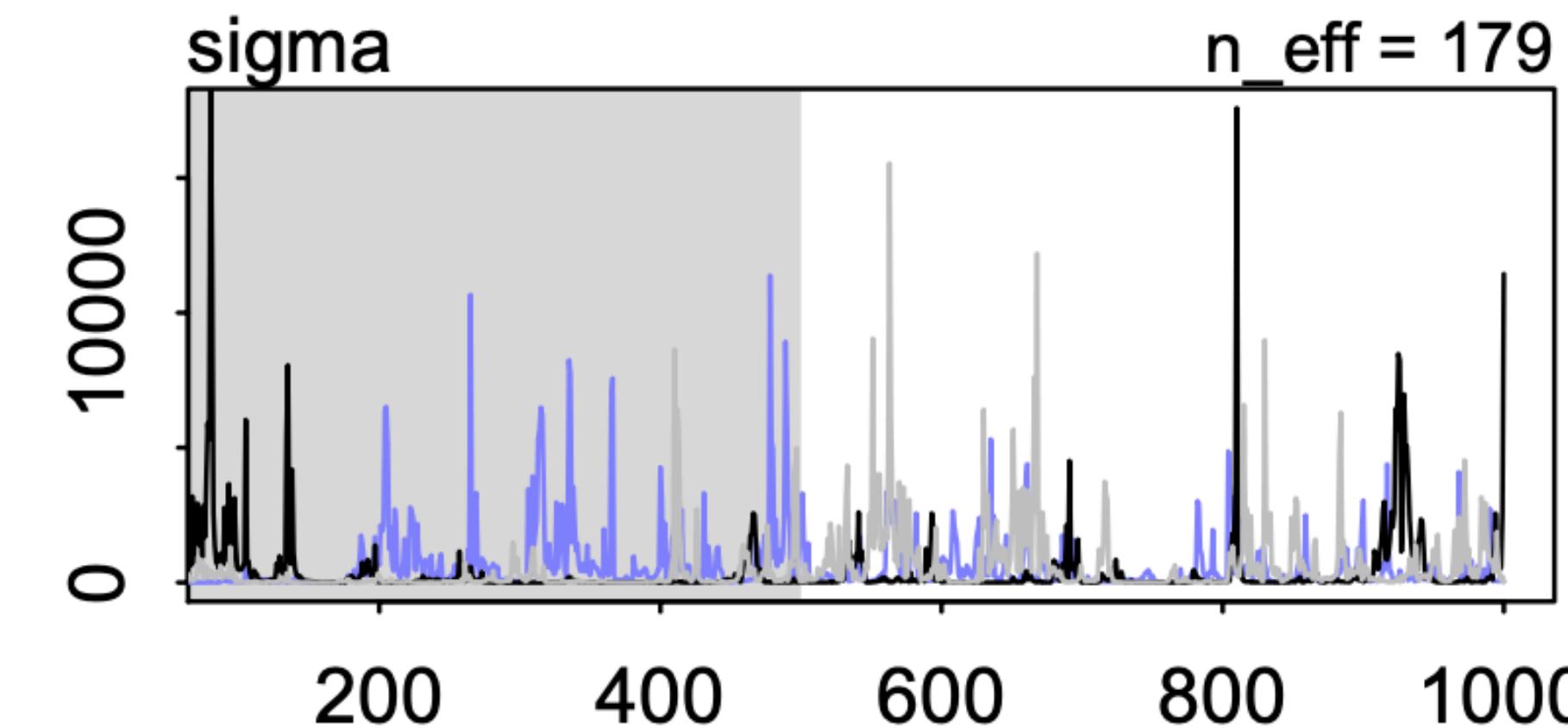
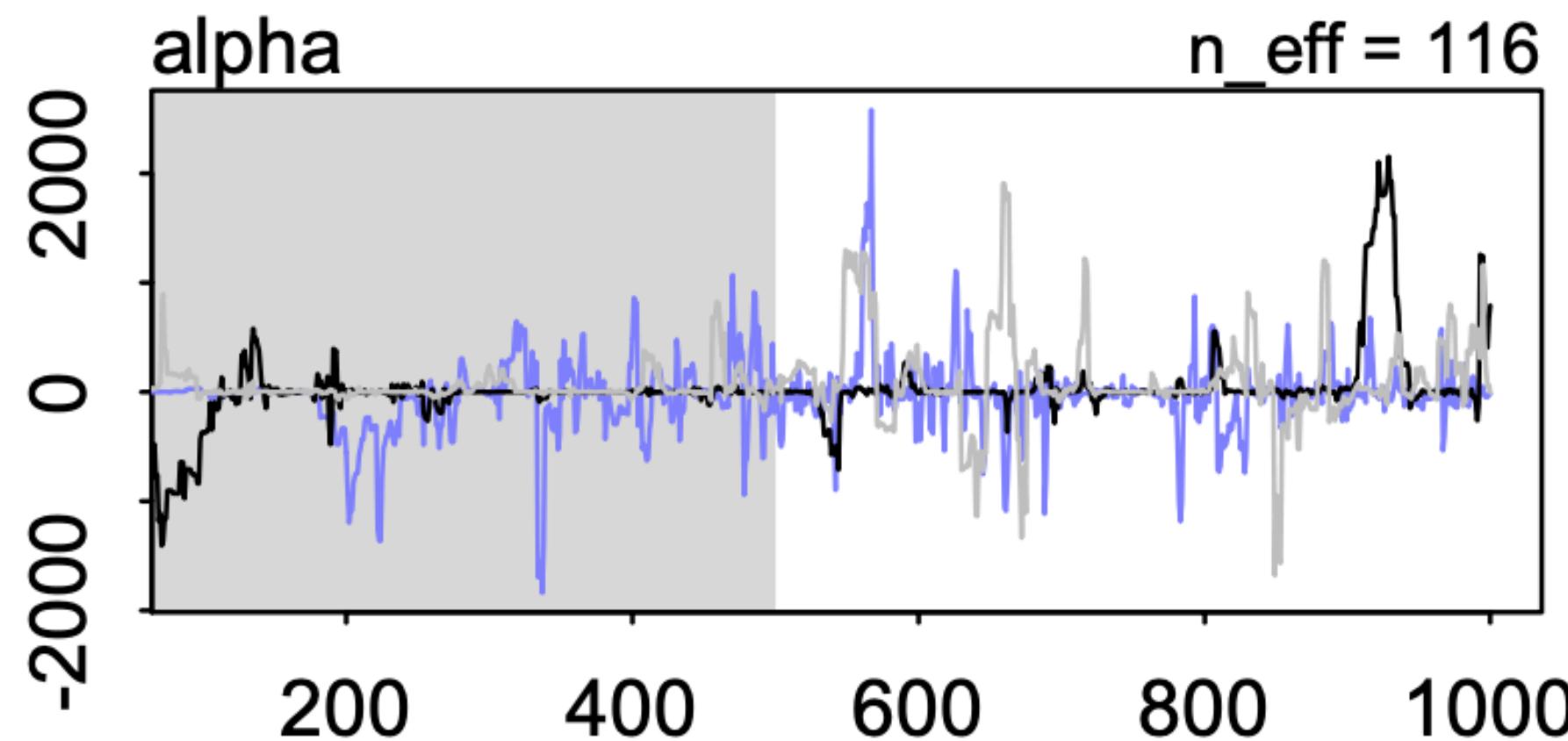
When samples are **autocorrelated**, you have fewer *effective* samples

```
> precis(mHMC)
    mean   sd  5.5% 94.5% n_eff Rhat4
a     0.00 0.10 -0.16  0.16  1632    1
bM    -0.06 0.17 -0.32  0.21  1137    1
bA    -0.61 0.17 -0.86 -0.34  1160    1
sigma 0.83 0.09  0.70  0.99  1504    1
>
```



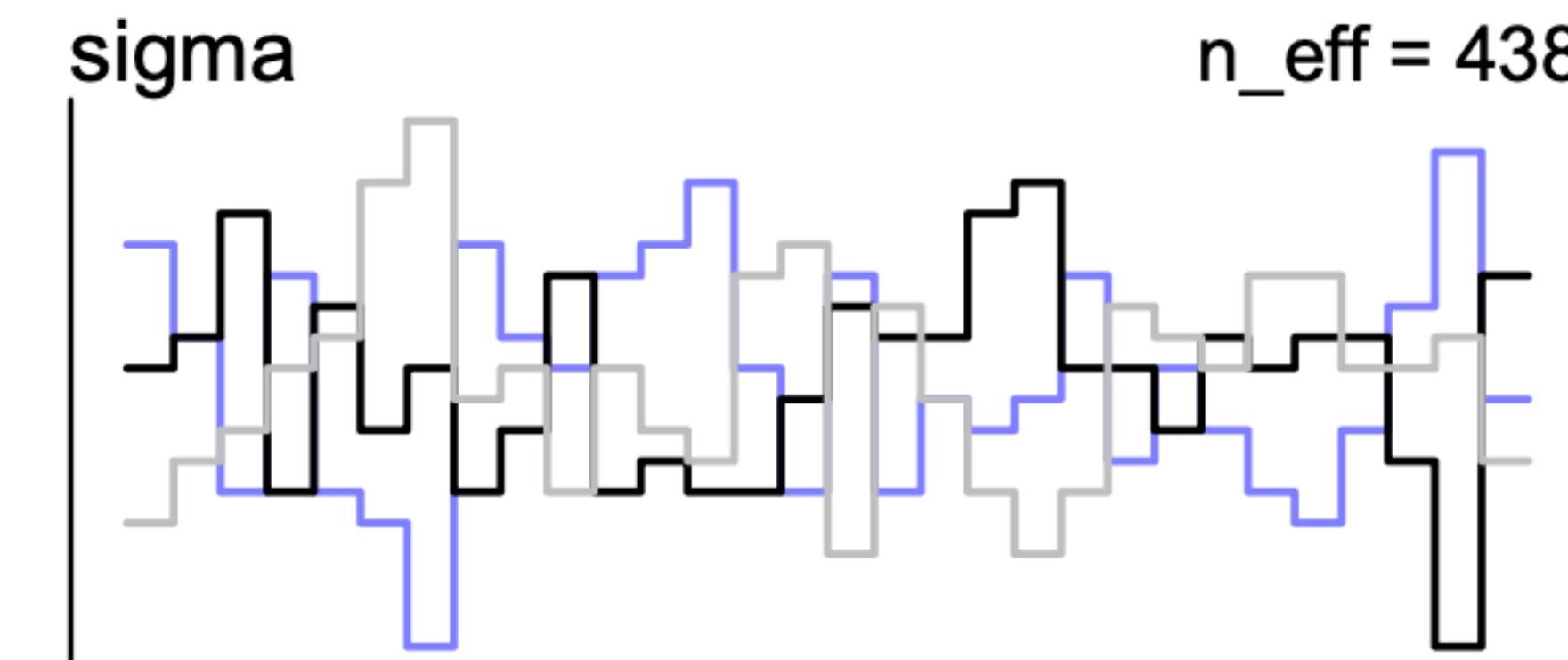
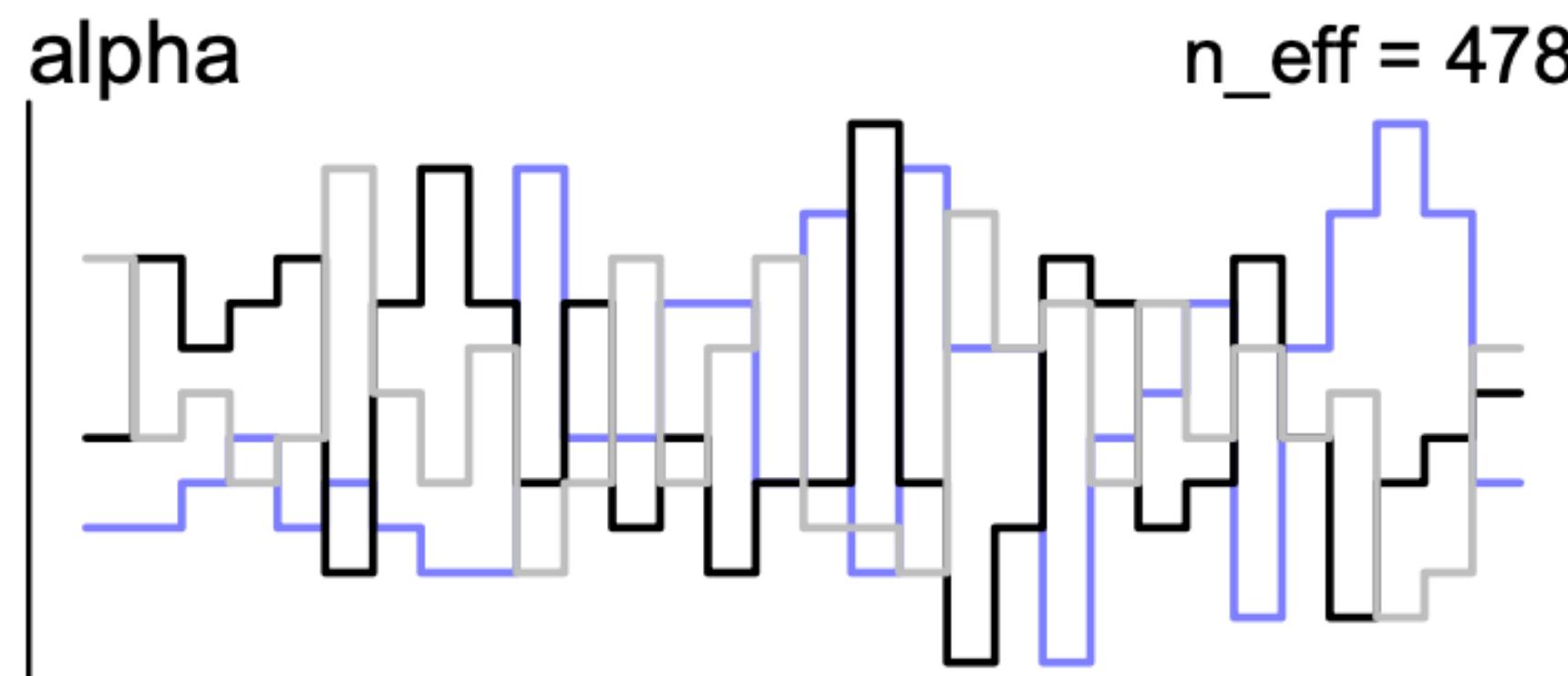
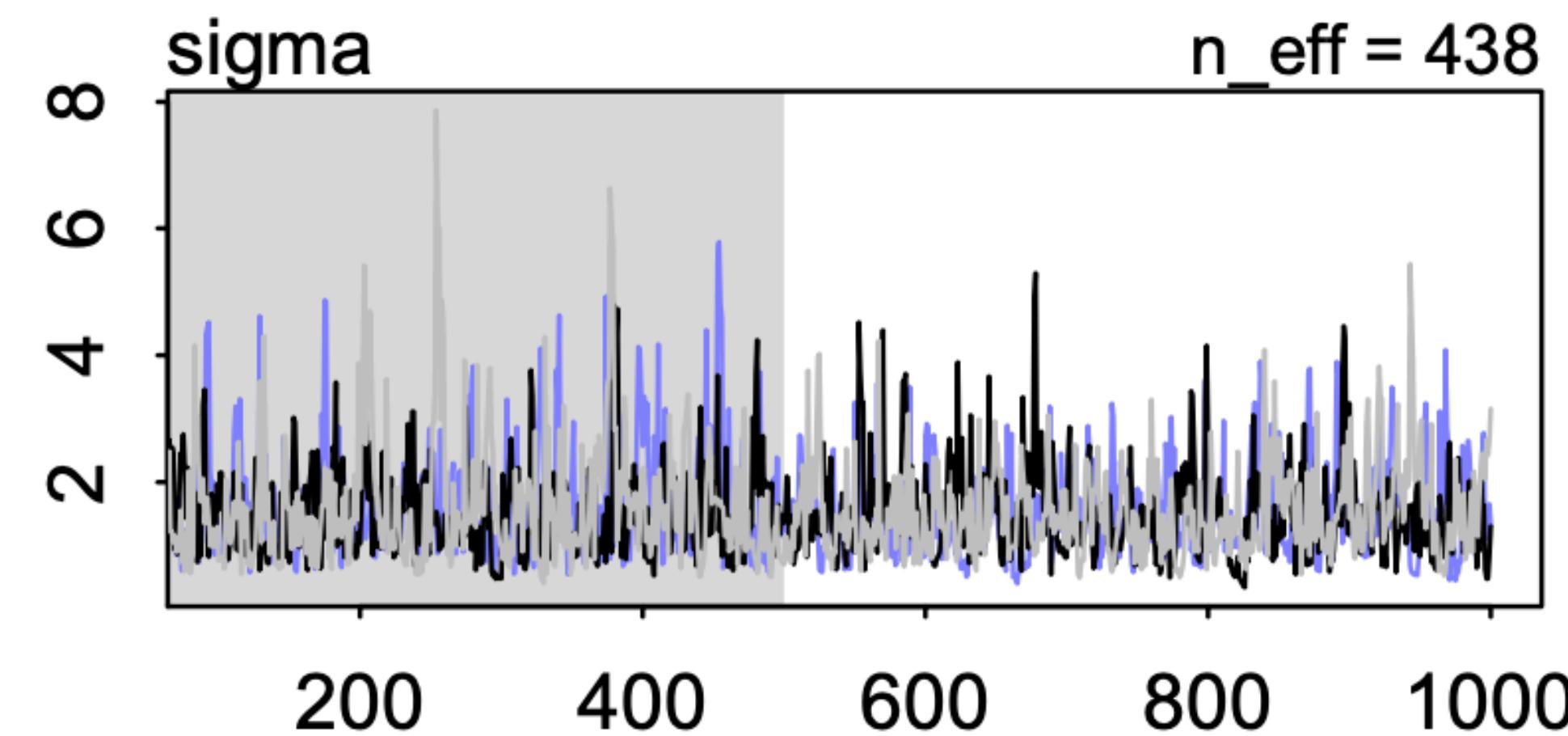
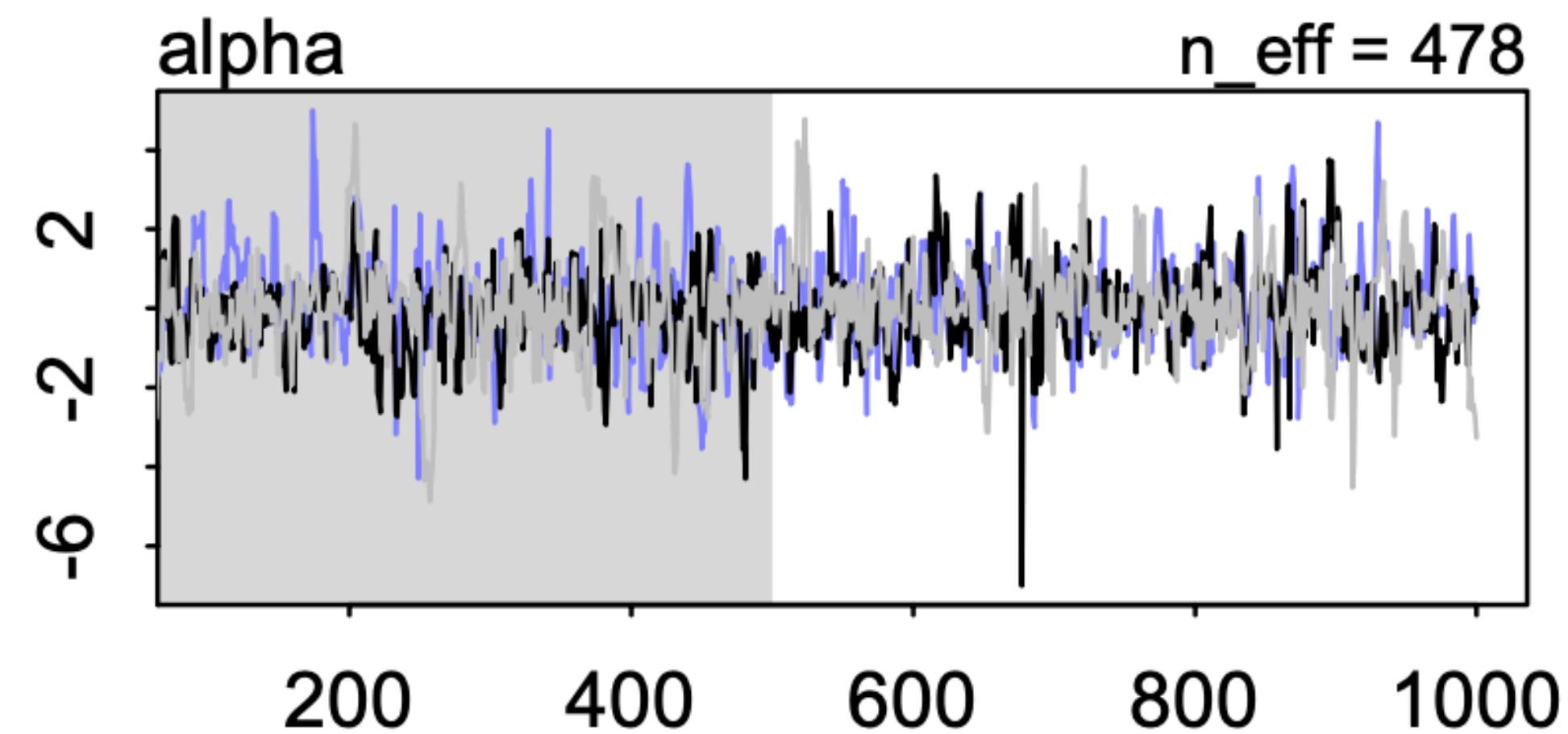
MCMC. Diagnosis of trace and trunk plots

- Not healthy ?



MCMC. Diagnosis of trace and trunk plots

- Healthy ?



Summary

- MCMC has been a revolution in scientific computing
- Custom scientific modeling
- High-dimension
- Do not “pipette by mouth”