See also 20 & 27 June 2018 at
http://phyloseminar.org/recorded.html

# Bayesian Phylogenetics

Workshop on Molecular Evolution
Woods Hole, Massachusetts

29 May 2022

## Paul O. Lewis
Department of Ecology & Evolutionary Biology
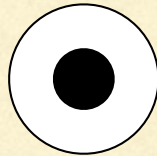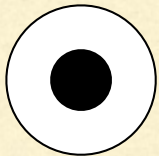
**UCONN**
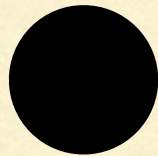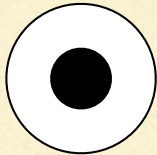UNIVERSITY OF CONNECTICUT

# Bayesian inference

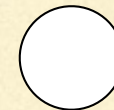# Joint probabilities

White,Solid

White,Dotted

Black,Dotted

Black,Solid

10 marbles in a bag
Sampling with replacement

$Pr(B,S) = 0.4$

$Pr(W,S) = 0.1$
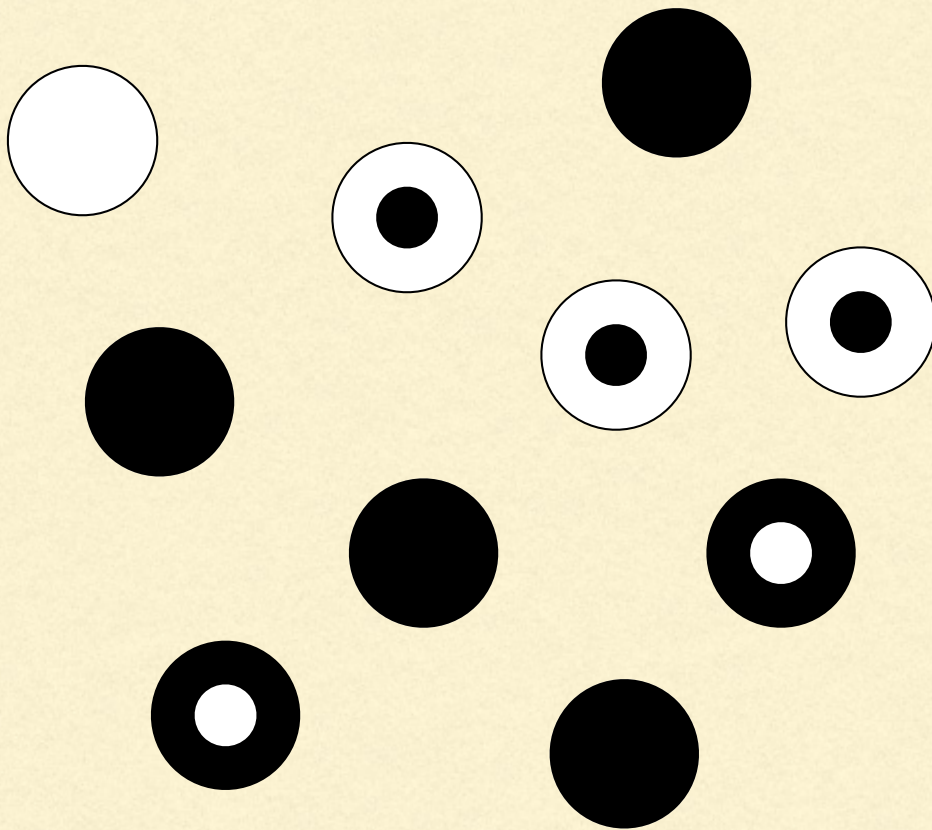
$Pr(B,D) = 0.2$

$Pr(W,D) = 0.3$

# Conditional probabilities

What's the probability that a marble is black given that it is dotted?

5 marbles satisfy the condition (D)

$$Pr(B|D) = \frac{2}{5}$$

2 remaining marbles are black (B)

# Marginal probabilities

W,D

B,D

Marginalizing over color yields the total probability that a marble is dotted (D)

$$Pr(\textbf{D}) = Pr(B,\textbf{D}) + Pr(W,\textbf{D})$$

$$= 0.2 + 0.3$$

$$= 0.5$$

Marginalization involves summing all joint probabilities containing D

# Marginalization

|   | B | W |
|---|---|---|
| D | Pr(D,B) | Pr(D,W) |
| S | Pr(S,B) | Pr(S,W) |

# Marginalizing over colors

# Joint probabilities

|   | B | W |
|---|---|---|
| D | Pr(D,B) | Pr(D,W) |
| S | Pr(S,B) | Pr(S,W) |

# Marginalizing over "dottedness"

# Bayes' rule



The joint probability Pr(B,D) can be written as the product of a *conditional probability* and the *probability of that condition*

$$Pr(B,D) \begin{cases} Pr(B|D)\ Pr(D) \\ Pr(D|B)\ Pr(B) \end{cases}$$

Either B or D can be the condition

# Bayes' rule

Equate the two ways of writing Pr(B,D)

$$Pr(B|D)\ Pr(D) = Pr(D|B)\ Pr(B)$$

Divide both sides by Pr(D)

$$\frac{Pr(B|D)\ \cancel{Pr(D)}}{\cancel{Pr(D)}} = \frac{Pr(D|B)\ Pr(B)}{Pr(D)}$$

Bayes' rule

$$Pr(B|D) = \frac{Pr(D|B)\ Pr(B)}{Pr(D)}$$

# Bayes' rule



$$\frac{2}{5} = \frac{\frac{1}{3} \times \frac{3}{5}}{\frac{1}{2}}$$

$$\frac{2}{5} = \frac{2}{5} \checkmark$$

Bayes' rule

$$Pr(B|D) = \frac{Pr(D|B) \, Pr(B)}{Pr(D)}$$

# Bayes' rule (variations)

$$\Pr(B|D) = \frac{\Pr(D|B)\Pr(B)}{\Pr(D)}$$

$$= \frac{\Pr(D|B)\Pr(B)}{\Pr(B,D) + \Pr(W,D)}$$

Pr(*D*) is the **marginal probability** of being dotted
To compute it, we **marginalize over colors**

# Bayes' rule (variations)

$$\Pr(B|D) = \frac{\Pr(D|B)\,\Pr(B)}{\Pr(B,D) + \Pr(W,D)}$$

$$= \frac{\Pr(D|B)\,\Pr(B)}{\Pr(D|B)\,\Pr(B) + \Pr(D|W)\,\Pr(W)}$$

$$= \frac{\Pr(D|B)\,\Pr(B)}{\sum_{\theta \in \{B,W\}} \Pr(D|\theta)\,\Pr(\theta)}$$

# Bayes' rule in statistics

**Likelihood** of hypothesis $\theta$

**Prior probability** of hypothesis $\theta$

$$\Pr(\theta|D) = \frac{\Pr(D|\theta)\,\Pr(\theta)}{\sum_{\theta}\Pr(D|\theta)\,\Pr(\theta)}$$

**Posterior probability** of hypothesis $\theta$

**Marginal probability of the data** (marginalizing over hypotheses)

# Paternity example

(father) A-     aa (mother)

Aa (child)

$$\Pr(\theta \mid D) = \frac{\Pr(D \mid \theta)\,\Pr(\theta)}{\sum_{\theta}\Pr(D \mid \theta)\,\Pr(\theta)}$$

| | $\theta_1$ | $\theta_2$ | Row sum |
|---|---|---|---|
| Genotypes | AA | Aa | --- |
| Prior | 1/2 | 1/2 | 1 |
| Likelihood | 1 | 1/2 | --- |
| Prior X Likelihood | 1/2 | 1/4 | 3/4 |
| Posterior | 2/3 | 1/3 | 1 |

# Bayes' rule: continuous case

Likelihood    Prior probability **density**

$$p(\theta \mid D) = \frac{p(D \mid \theta)\,p(\theta)}{\int p(D \mid \theta)\,p(\theta)\,d\theta}$$

Posterior probability
**density**

Marginal probability
of the data

# If you had to guess...



Photo by Tracy Heath

**1 meter**

*Not knowing anything about my archery abilities*, draw a curve representing your view of the chances of my arrow landing a distance *d* centimeters from the center of the target.

0.0

$d$ (centimeters from target center)

# Case 1: assume I have talent

An *informative* prior (low variance) that says most of my arrows will fall within 20 cm of the center (thanks for your confidence!)

1 meter

0.0          20.0          40.0          60.0          80.0

Case 2: assume I have a talent for missing the target!

1 meter

Also an *informative* prior, but one that says most of my arrows will fall within a narrow range just outside the entire target!

0.0          20.0          40.0          60.0

# Case 3: assume I have no talent

This is a *vague* prior: its **high variance** reflects nearly total ignorance of my abilities, saying that my arrows could land nearly anywhere!

1 meter

0.0  20.0  40.0  60.0  80.0

# A matter of scale

Notice that I haven't provided a scale for the vertical axis.

What exactly does the height of this curve mean?

For example, does the height of the dotted line represent the *probability* that my arrow lands 60 cm from the center of the target?

∞

0.0          20.0          40.0          60.0

# Probabilities are associated with intervals

**Probabilities** are attached to **intervals** (i.e. ranges of values), **not** individual **values**

The probability of any given point (e.g. $d = 60.0$) is zero!

However, we can ask about the probability that $d$ falls in a particular interval e.g. $50.0 < d < 65.0$

0.0                    20.0                    40.0                    60.0

Probabilities vs. probability densities

Probability **density** function

Note: the height of this curve does **not** represent a probability (if it did, it would not exceed 1.0)

# Densities of various substances

| Substance | Density (g/cm$^3$) |
|---|---|
| Cork | 0.24 |
| Aluminum | 2.7 |
| Gold | 19.3 |

*Density does not equal mass*
mass = density × volume

# A brick with varying density

density

distance from left end

Gold on
left end

Aluminum on
right end

# Integrating a density yields a probability

Area of rectangle = $p(\theta)d\theta$

$p(\theta)$

$d\theta$

infinitesimal width

Long s from U.S. Bill of Rights

$$1.0 = \int p(\theta)d\theta$$

The density curve is scaled so that the value of this integral (i.e. the total area) equals 1.0

$p(\theta)$

$\theta$

# Integrating a density yields a probability

Area of rectangle $= p(\theta)d\theta$

$p(\theta)$

$d\theta$

infinitesimal
width

$p(\theta)$

$\theta$

$$0.39109 = \int_1^2 p(\theta)d\theta$$

The **area** under the density curve from 1 to 2 is the **probability** that $\theta$ is between 1 and 2

# Archery priors revisited



mean=1.732
variance=3

mean=60
variance=3

mean=200
variance=40000

These density curves are all variations of a **gamma probability distribution**. We could have used a gamma distribution to specify each of the prior probability distributions for the archery example. Note that **higher variance** means **less informative**

# Usually there are many parameters...
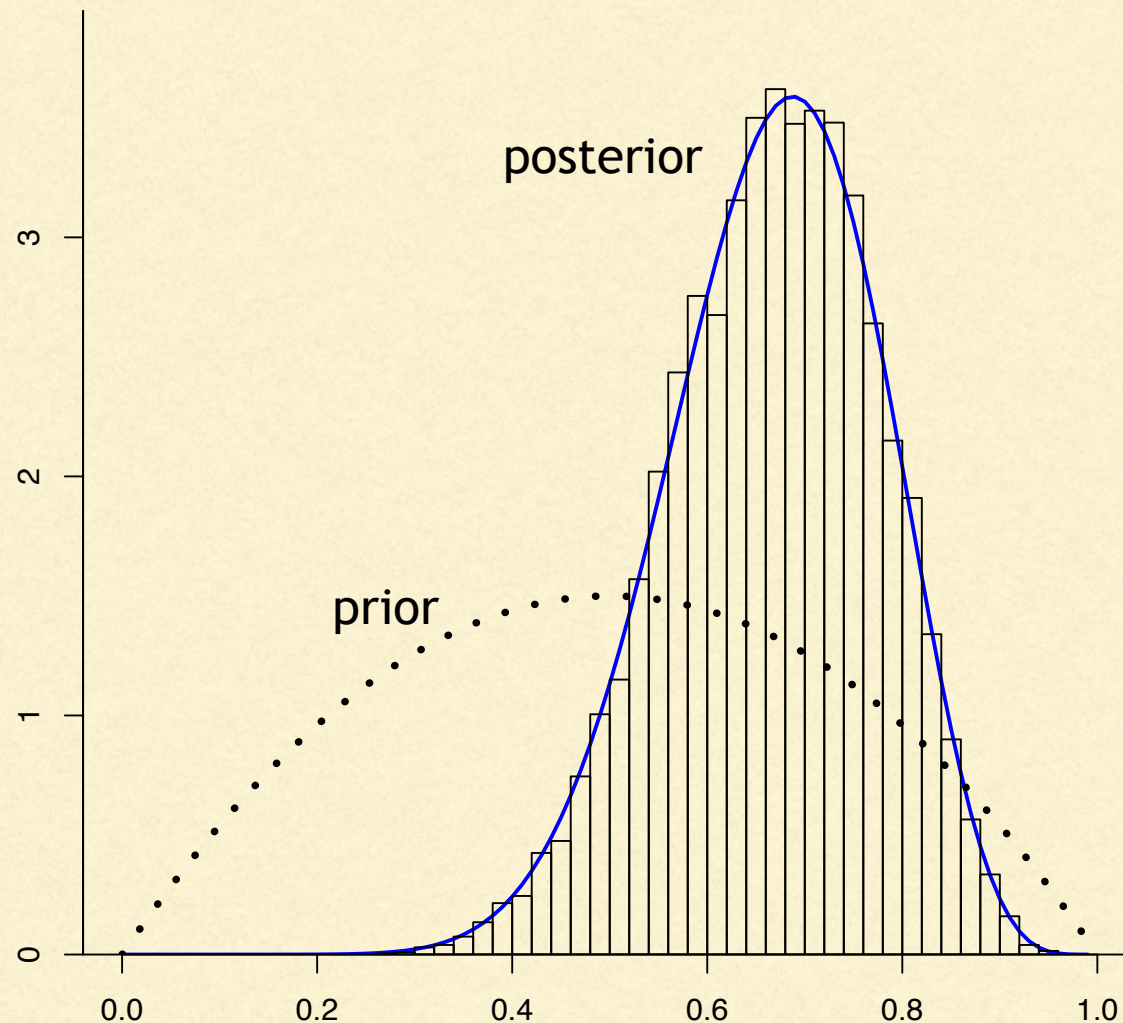
A 2-parameter example

Likelihood   Prior density

$$p(\theta, \phi \,|\, D) = \frac{p(D\,|\,\theta, \phi)\; p(\theta)\; p(\phi)}{\int_\theta \int_\phi p(D\,|\,\theta, \phi)\; p(\theta)\; p(\phi)\; d\phi\; d\theta}$$

Posterior probability density

Marginal probability of data

An analysis of **100 sequences** under the simplest model (JC69) requires 197 branch length parameters. The denominator would require a **197-fold integral** inside a sum over **all possible tree topologies**! It would thus be nice to avoid having to calculate the marginal probability of the data...

# Markov chain Monte Carlo (MCMC)

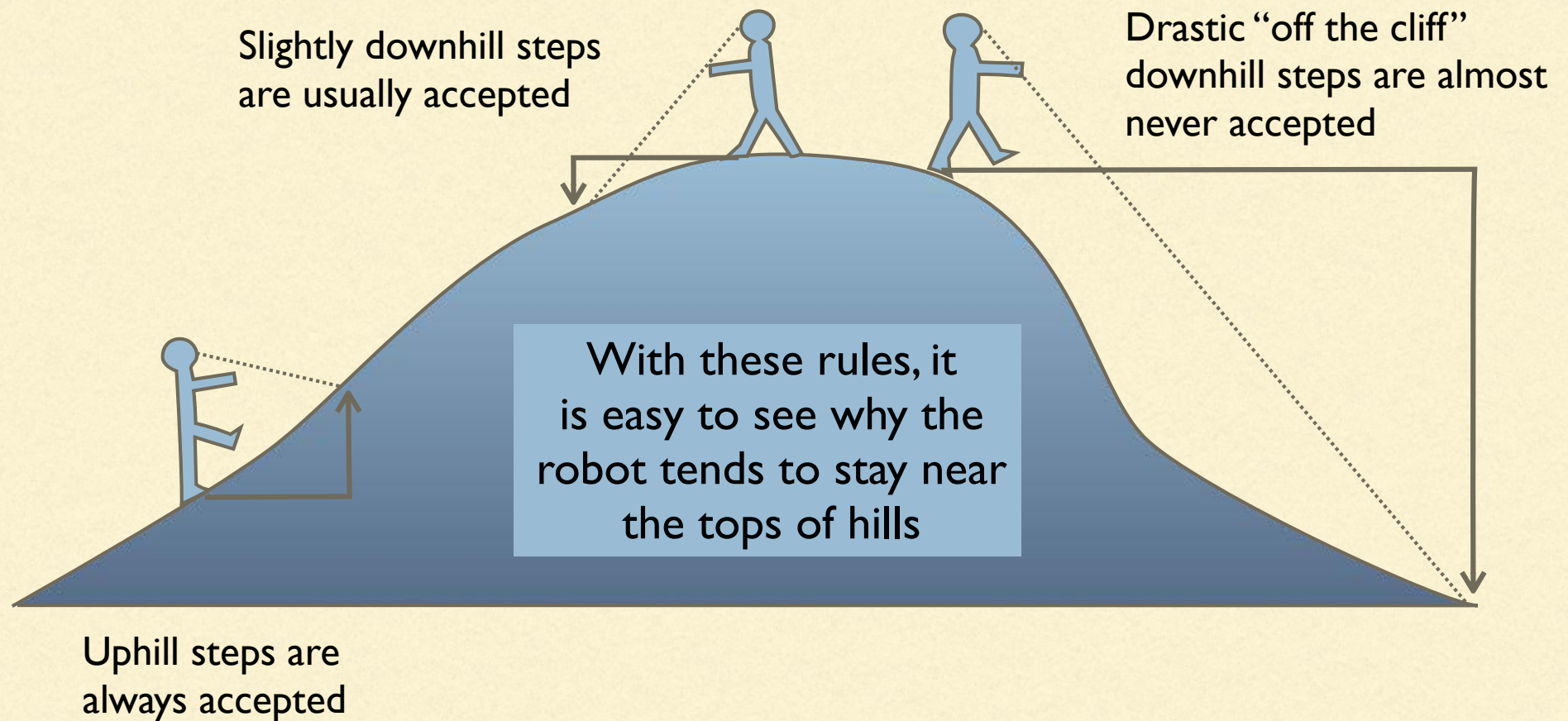# Markov chain Monte Carlo (MCMC)
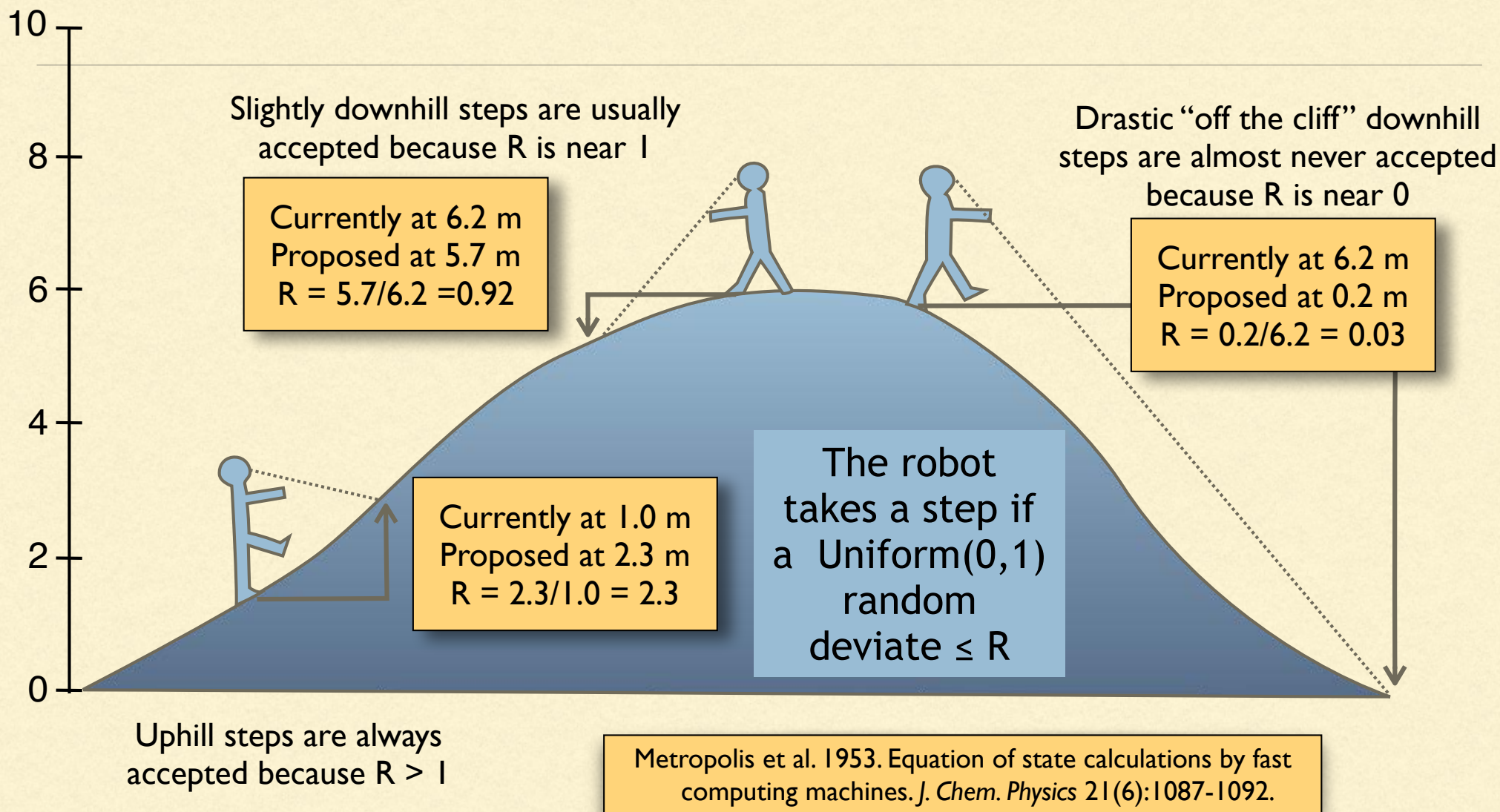


For more complex problems, we might settle for a

## good approximation

to the posterior distribution

# MCMC robot's rules



Slightly downhill steps are usually accepted

Drastic "off the cliff" downhill steps are almost never accepted

With these rules, it is easy to see why the robot tends to stay near the tops of hills

Uphill steps are always accepted

# Actual rules (Metropolis algorithm)

Slightly downhill steps are usually accepted because R is near 1

Currently at 6.2 m
Proposed at 5.7 m
R = 5.7/6.2 =0.92

Drastic "off the cliff" downhill steps are almost never accepted because R is near 0

Currently at 6.2 m
Proposed at 0.2 m
R = 0.2/6.2 = 0.03

Currently at 1.0 m
Proposed at 2.3 m
R = 2.3/1.0 = 2.3

The robot takes a step if a Uniform(0,1) random deviate ≤ R

Uphill steps are always accepted because R > 1

Metropolis et al. 1953. Equation of state calculations by fast computing machines. *J. Chem. Physics* 21(6):1087-1092.

# Cancellation of marginal likelihood

When calculating the ratio ($R$) of posterior densities, the marginal probability of the data cancels.

$$\frac{p(\theta^*|D)}{p(\theta|D)} = \frac{\dfrac{p(D|\theta^*)\,p(\theta^*)}{p(D)}}{\dfrac{p(D|\theta)\,p(\theta)}{p(D)}} = \frac{p(D|\theta^*)\,p(\theta^*)}{p(D|\theta)\,p(\theta)}$$

**Posterior odds**

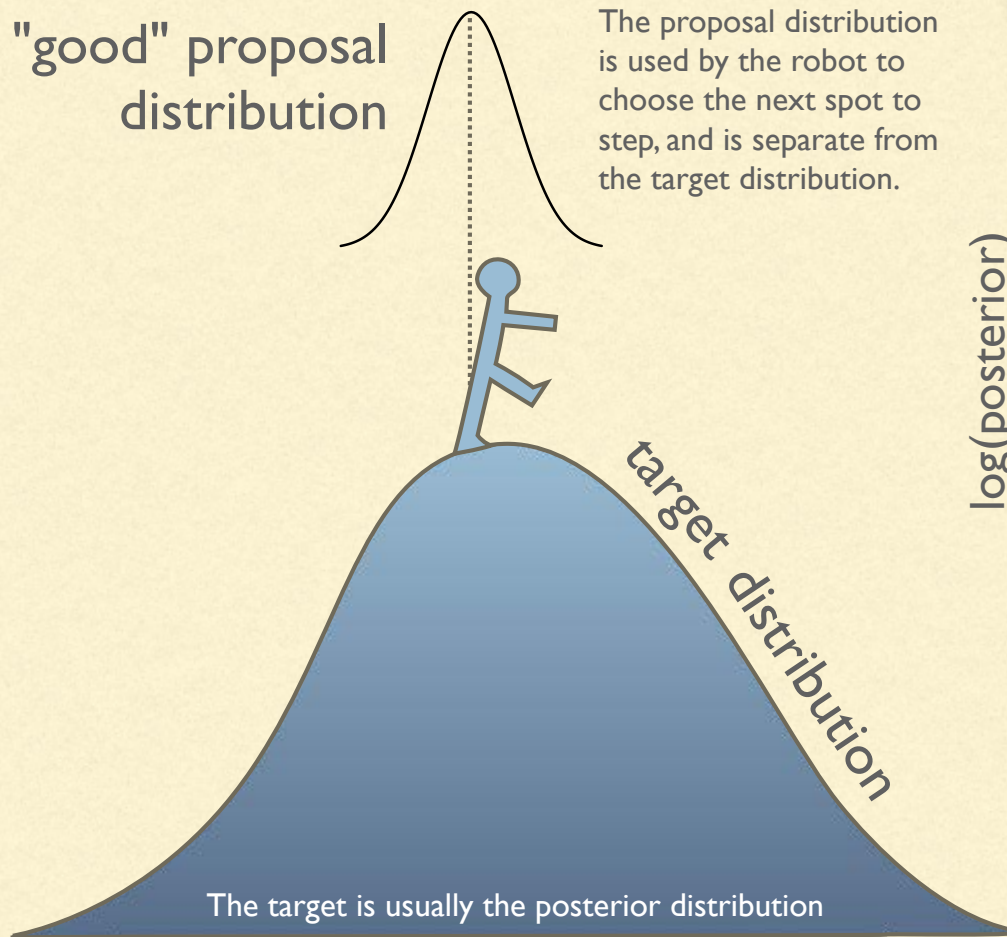Apply Bayes' rule to both top and bottom
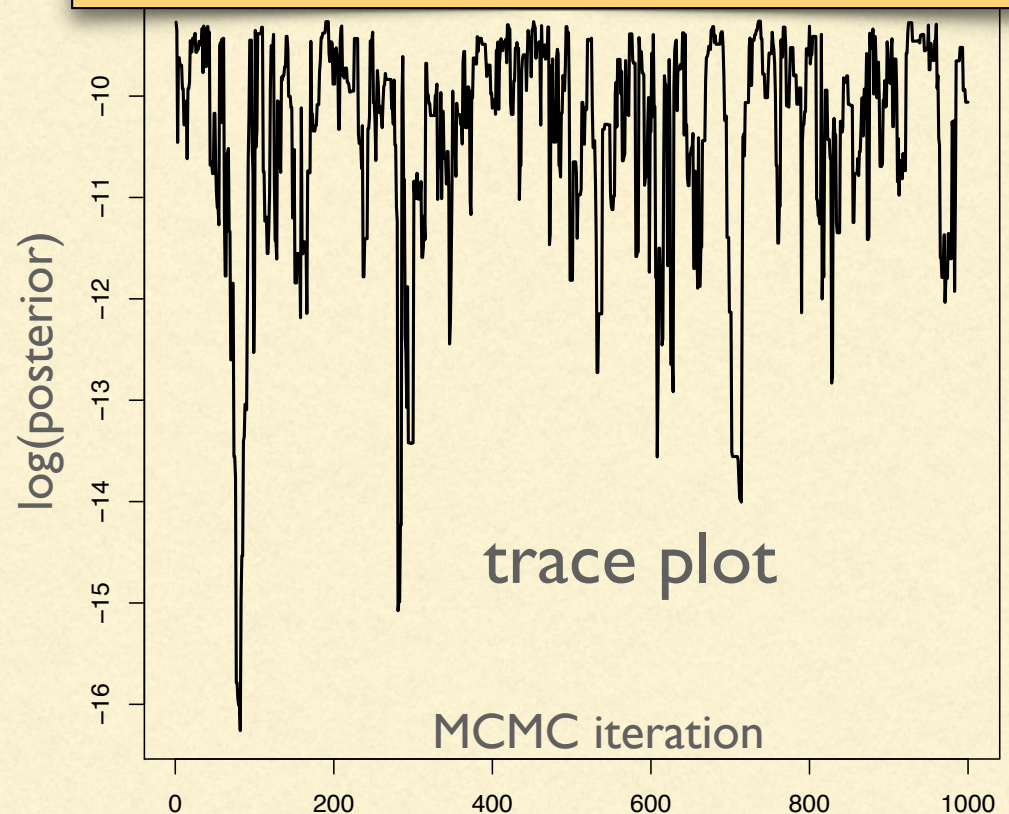
**Likelihood ratio**

**Prior odds**

# Target vs. Proposal Distributions

"good" proposal distribution

The proposal distribution is used by the robot to choose the next spot to step, and is separate from the target distribution.
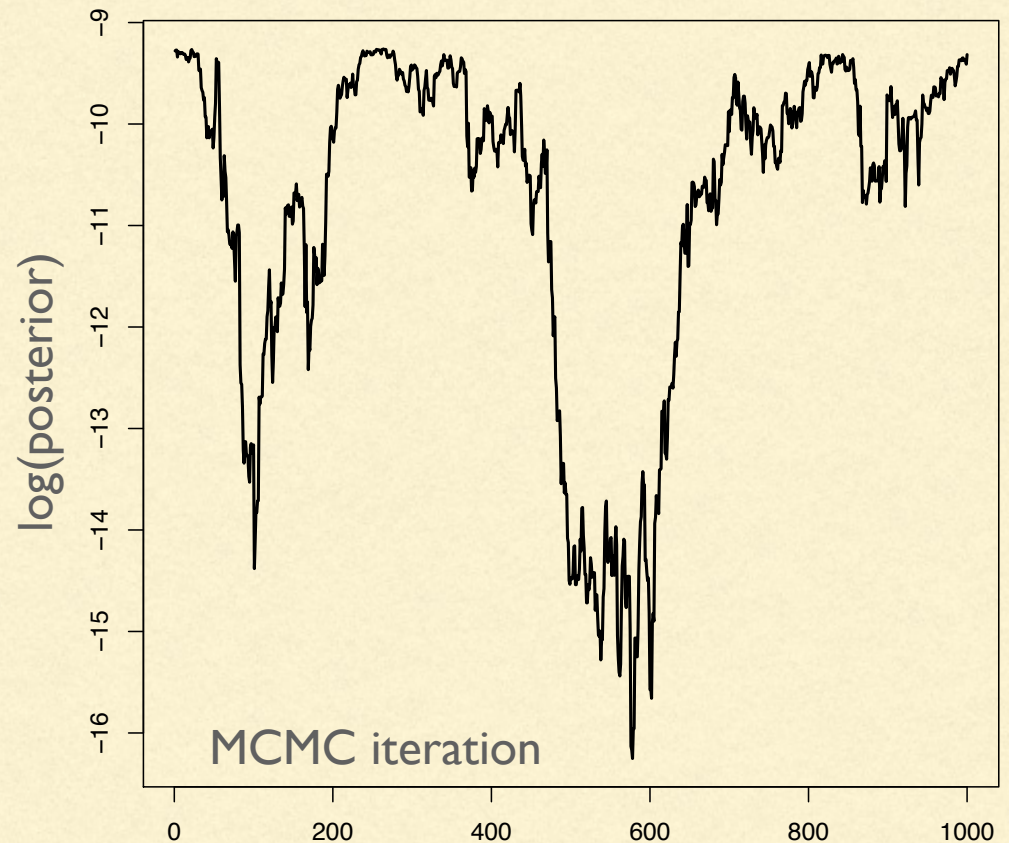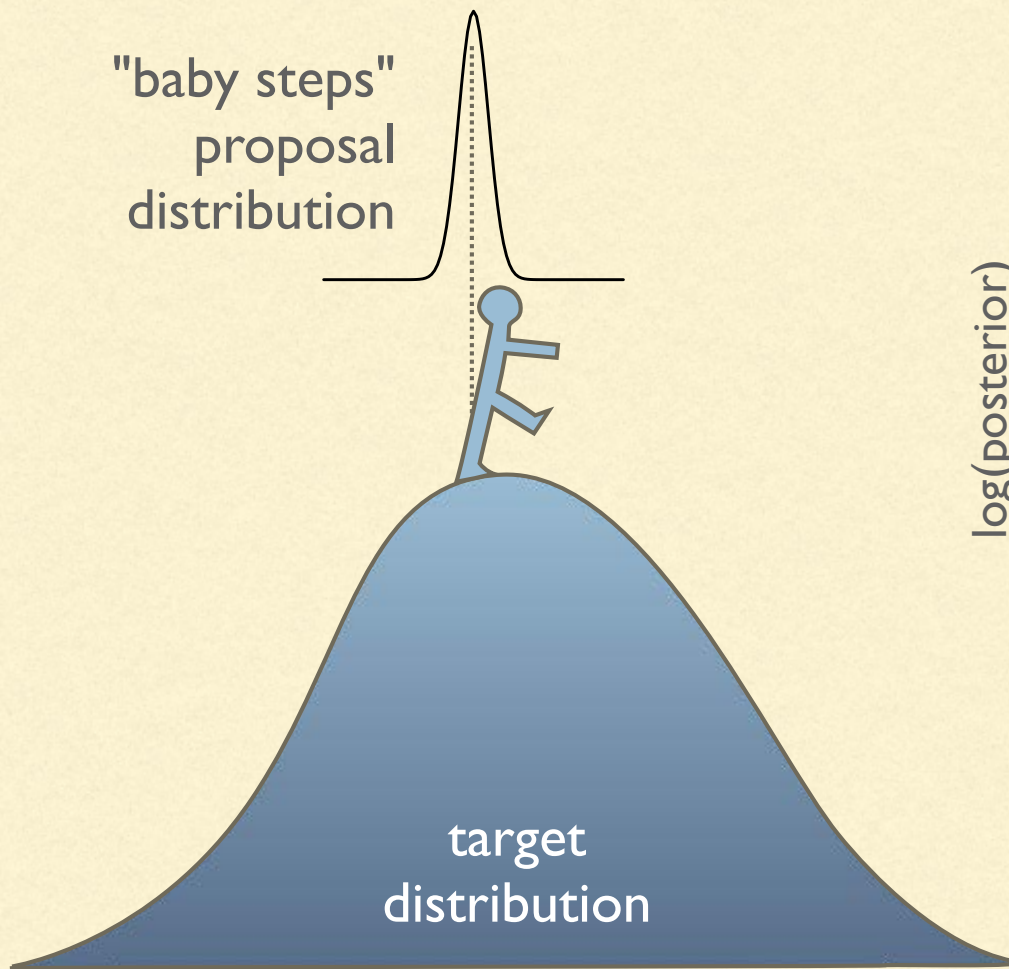
target distribution

The target is usually the posterior distribution

trace plot

log(posterior)

MCMC iteration

White noise appearance is a sign of good mixing

# Target vs. Proposal Distributions

"baby steps"
proposal
distribution

target
distribution



log(posterior)

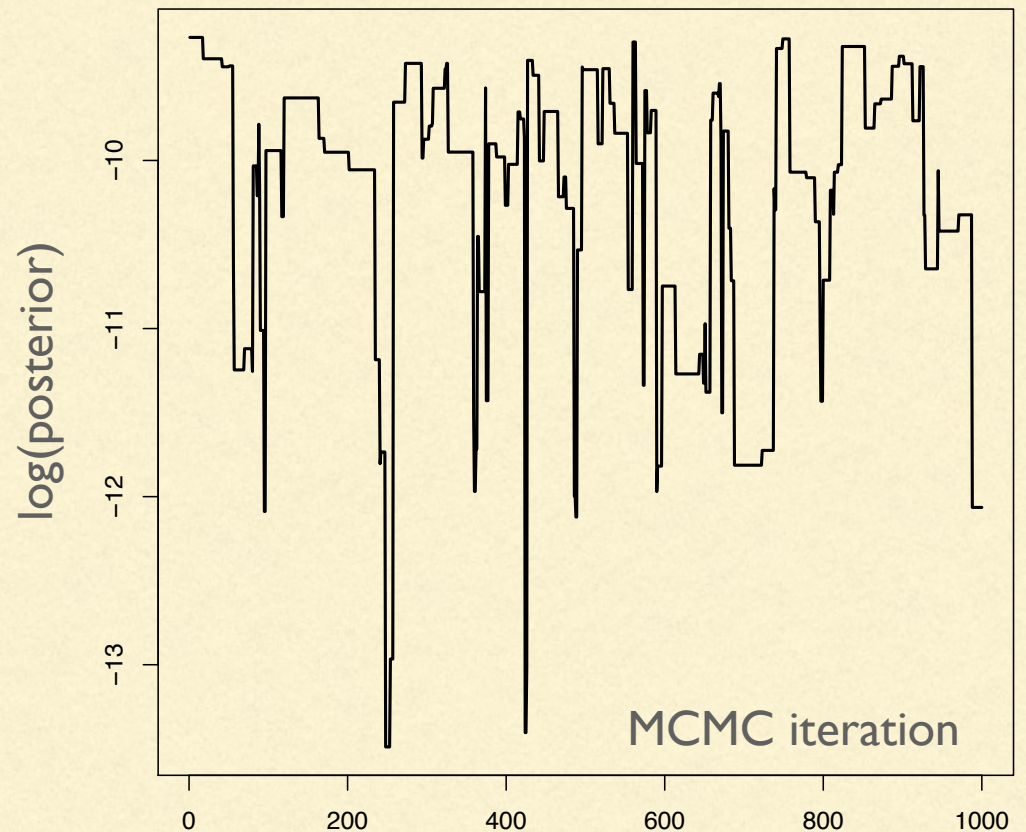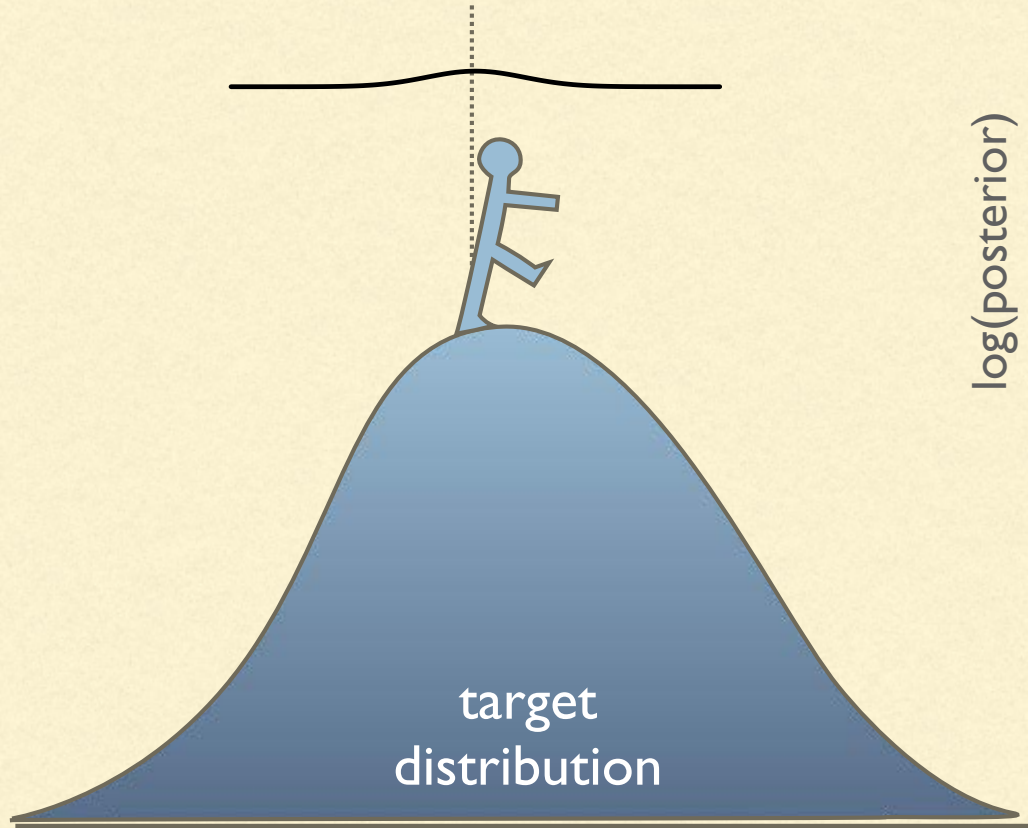MCMC iteration

Big waves in trace plot indicate
robot is crawling around

# Target vs. Proposal Distributions



"overly bold" proposal distribution

target distribution

log(posterior)

−10   −11   −12   −13

MCMC iteration

0    200    400    600    800    1000
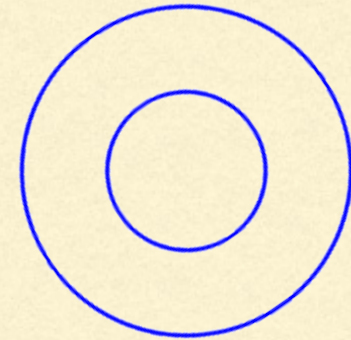
Plateaus in trace plot indicate robot is often stuck in one place

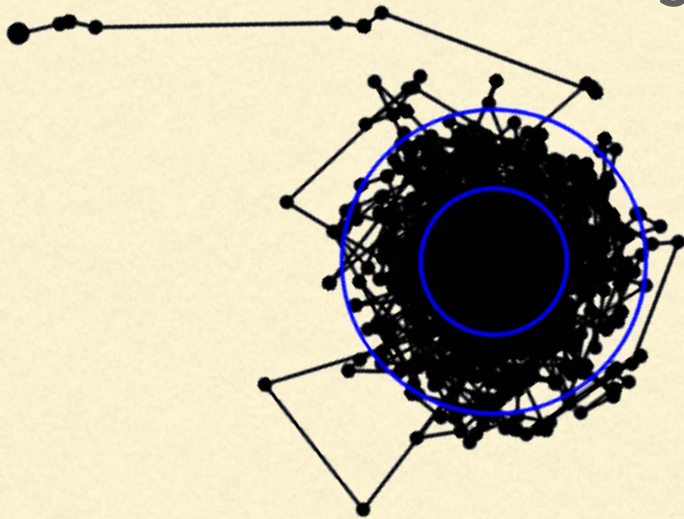# MCRobot (or "MCMC Robot")

Javascript version used today will run in most web
browsers and is available here:

https://plewis.github.io/applets/mcmc-robot/

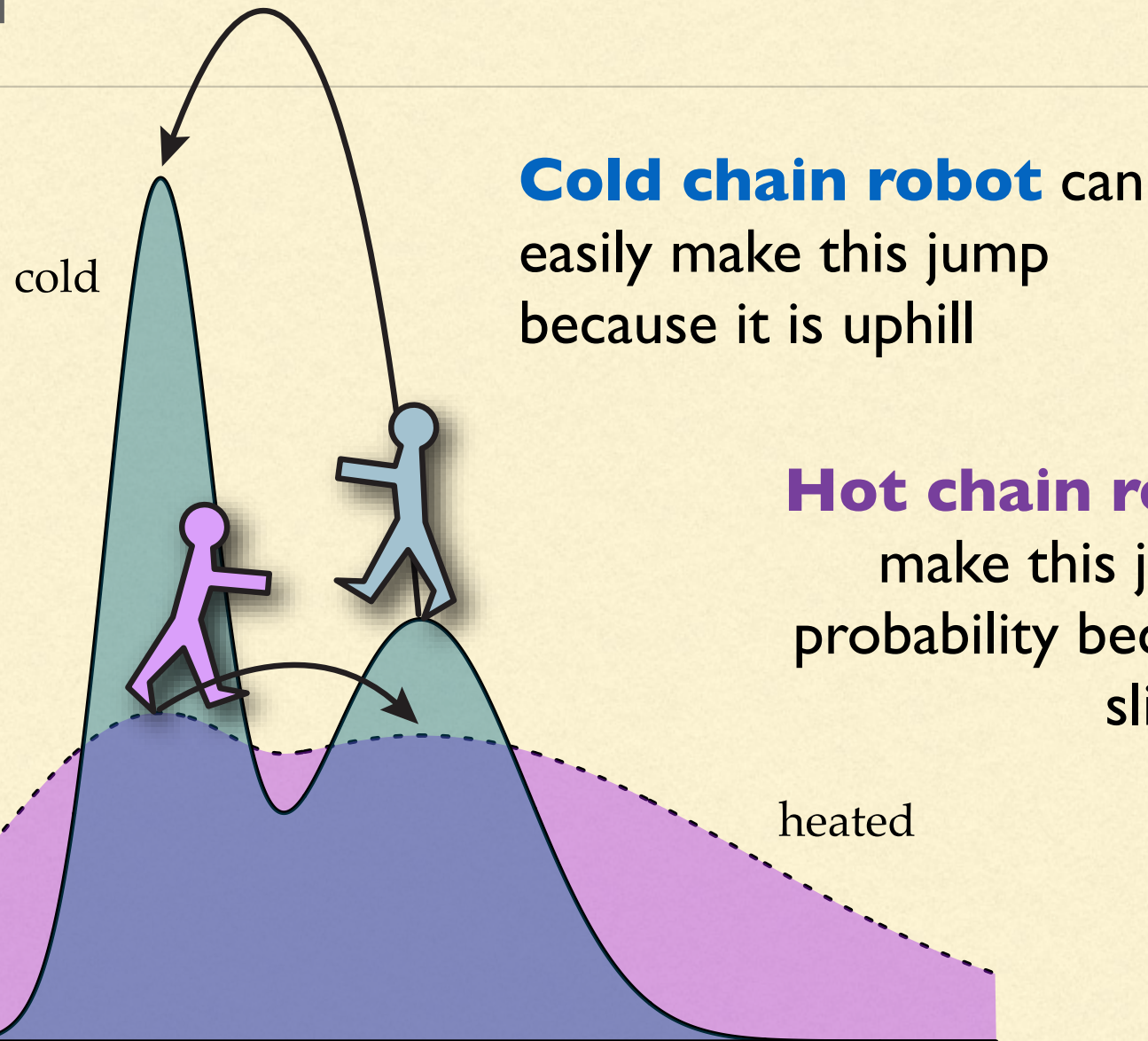# Metropolis-coupled Markov chain Monte Carlo (MCMCMC)

Sometimes the robot needs some help,

MCMCMC introduces helpers in the form of "heated chain" robots that can act as scouts.
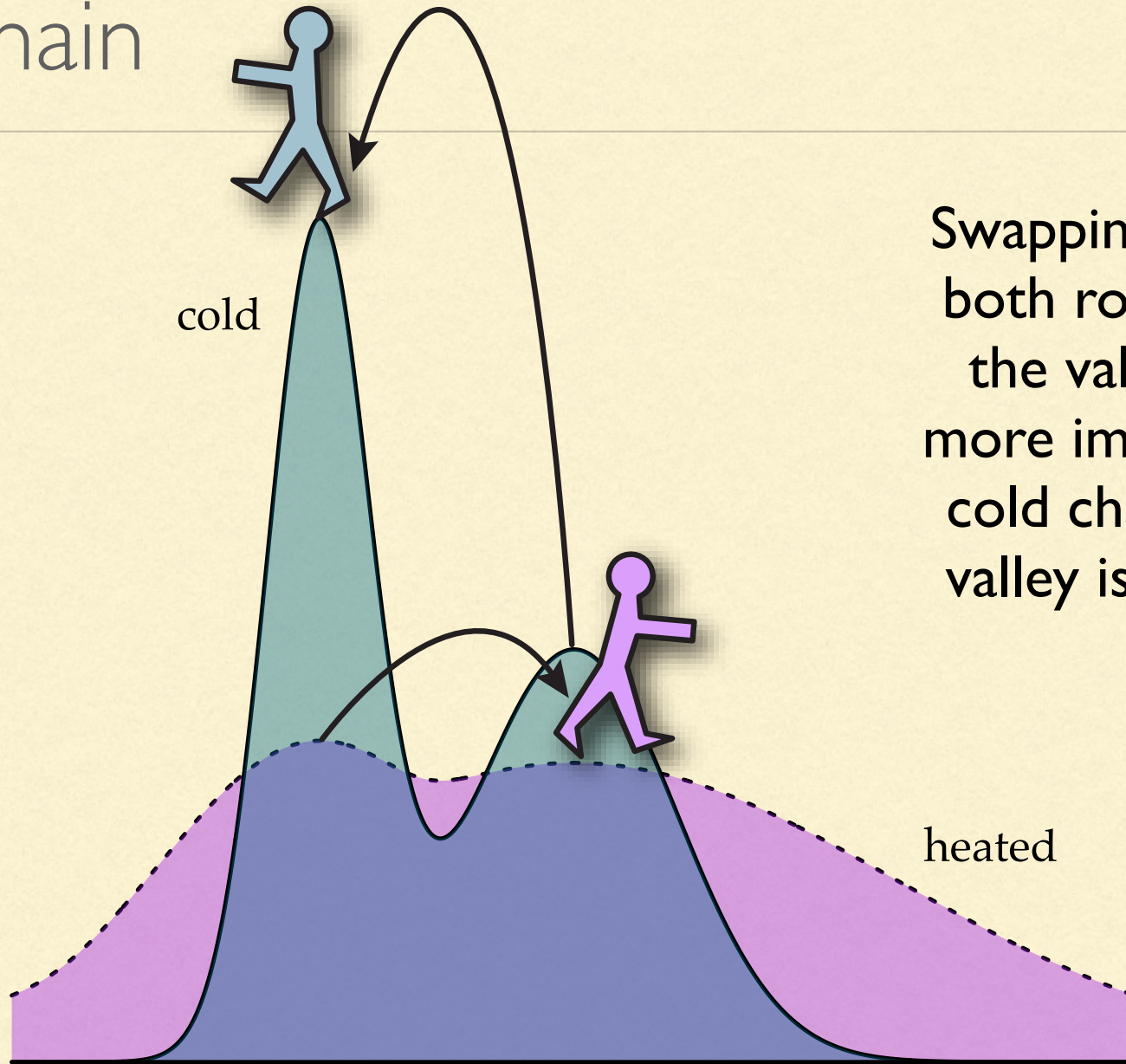
Geyer, C. J. 1991. Markov chain Monte Carlo maximum likelihood for dependent data. Pages 156-163 *in* Computing Science and Statistics (E. Keramidas, ed.).

# Heated chains act as scouts for the cold chain

cold

**Cold chain robot** can easily make this jump because it is uphill
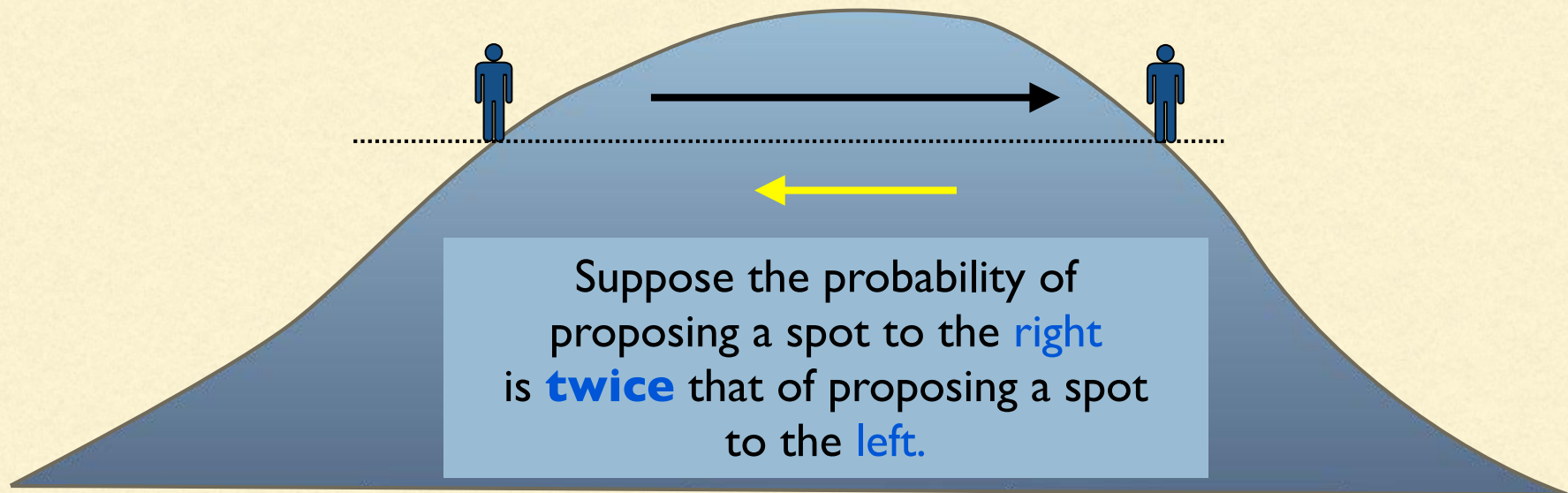
**Hot chain robot** can also make this jump with high probability because it is only slightly downhill
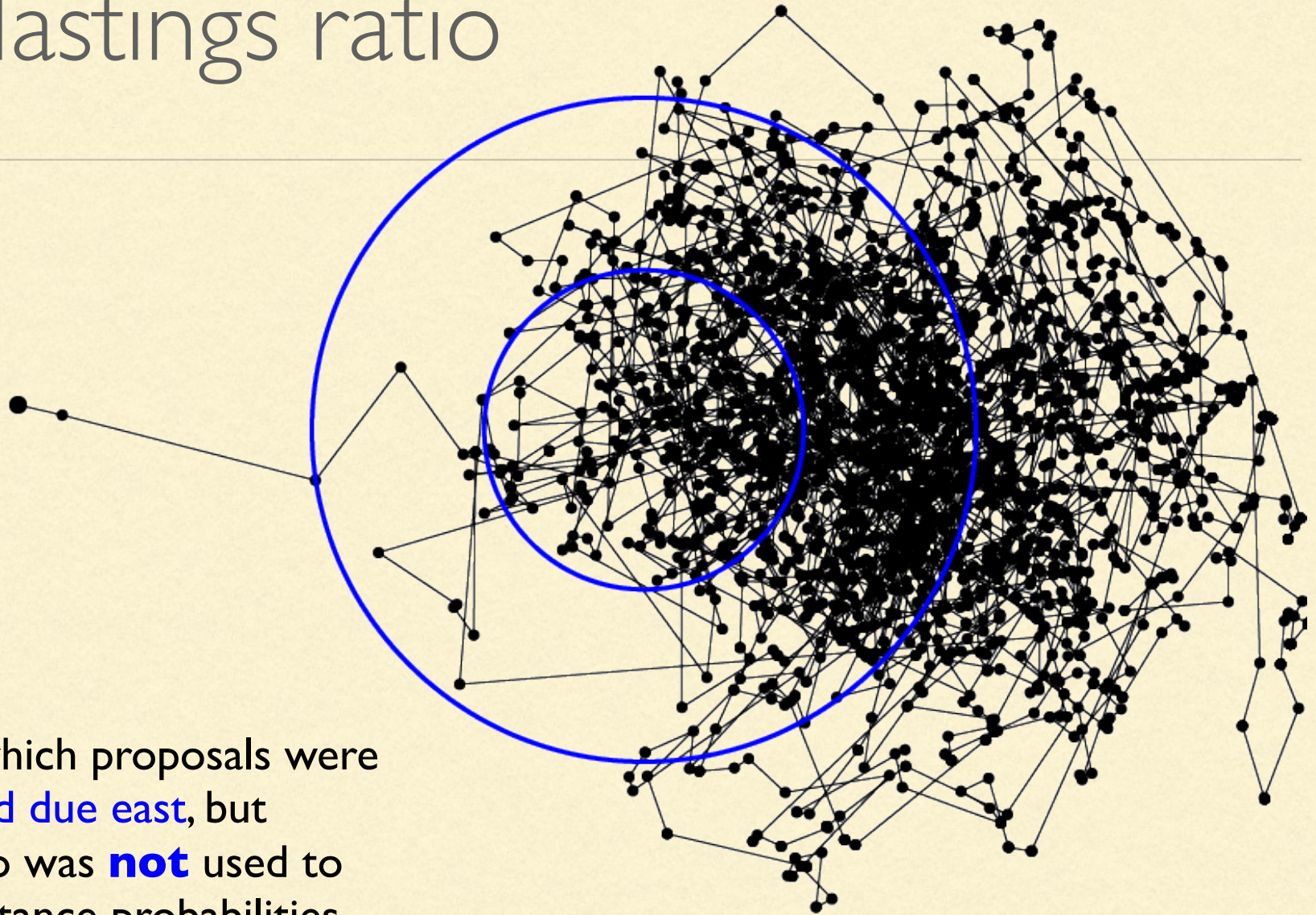
heated

# Heated chains act as scouts for the cold chain

cold

Swapping places means both robots can cross the valley, but this is more important for the cold chain because its valley is much deeper.

heated

MOLE 2022 Paul O. Lewis

# The Hastings ratio



Suppose the probability of proposing a spot to the right is **twice** that of proposing a spot to the left.
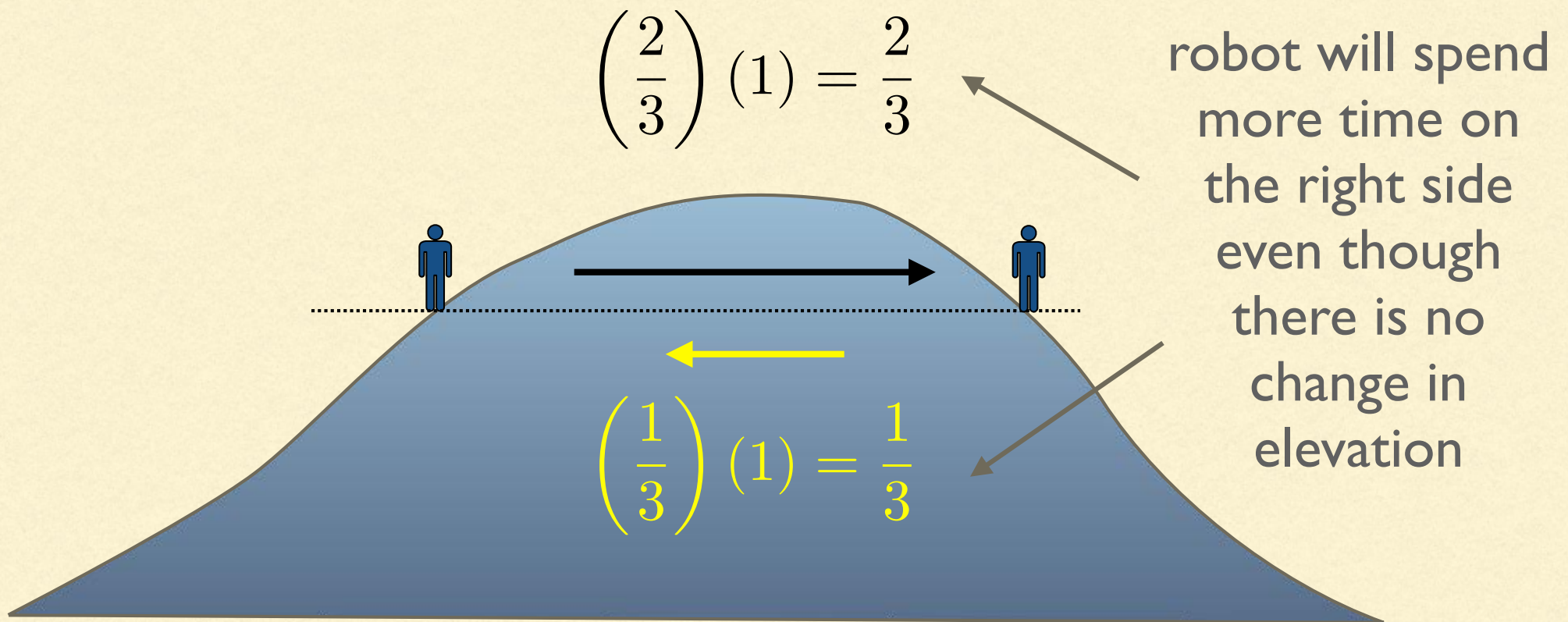
Hastings, W. K. 1970. Monte Carlo sampling methods using Markov chains and their applications. Biometrika 57:97-109.

# The Hastings ratio



Example in which proposals were
biased toward due east, but
Hastings ratio was **not** used to
modify acceptance probabilities

# The Hastings ratio

$$\left(\frac{2}{3}\right)(1) = \frac{2}{3}$$

$$\left(\frac{1}{3}\right)(1) = \frac{1}{3}$$

robot will spend more time on the right side even though there is no change in elevation

Hastings, W. K. 1970. Monte Carlo sampling methods using Markov chains and their applications. Biometrika 57:97-109.

# The Hastings ratio

$$\left(\frac{2}{3}\right)\left[(1)\left(\frac{1/3}{2/3}\right)\right] = \frac{1}{3}$$

robot spends same amount of time on both sides, as it should

$$\left(\frac{1}{3}\right)\left[(1)\left(\frac{2/3}{1/3}\right)\right] = \frac{1}{3}$$

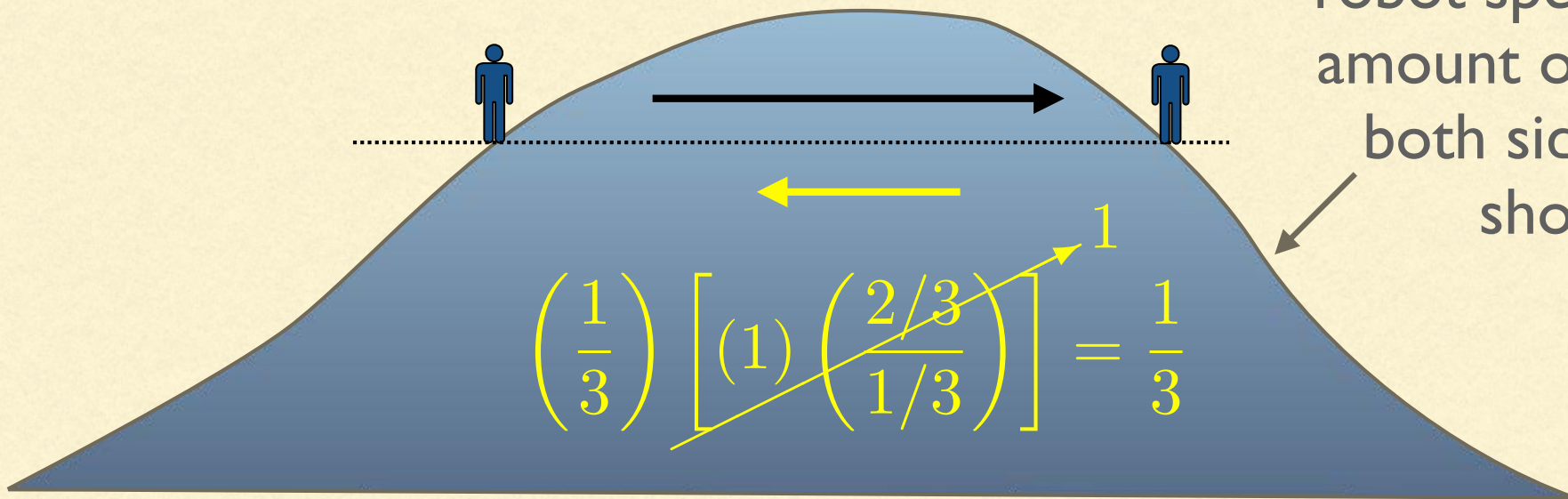Hastings, W. K. 1970. Monte Carlo sampling methods using Markov chains and their applications. Biometrika 57:97-109.

# Hastings Ratio

$$R = \min \left\{ 1, \left[ \frac{p(D|\theta^*)\, p(\theta^*)}{p(D|\theta)\, p(\theta)} \right] \left[ \frac{q(\theta|\theta^*)}{q(\theta^*|\theta)} \right] \right\}$$

$$\underbrace{\qquad\qquad\qquad}_{\text{posterior ratio}} \quad \underbrace{\qquad\qquad}_{\text{Hastings ratio}}$$

Note that the Hastings ratio is 1.0 if $q(\theta^*|\theta) = q(\theta|\theta^*)$