
Phylogenetics 101



Part 3a: Bayesian statistics and MCMC

- Bayesian inference
- Markov chain Monte Carlo
- MCMC in phylogenetics

Paul O. Lewis
Dept. Ecol. & Evol. Biology
University of Connecticut
<https://phylogeny.uconn.edu>

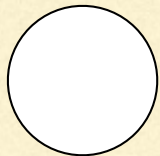


DEB-1354146 Broader Impacts

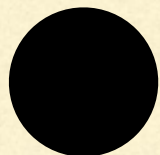
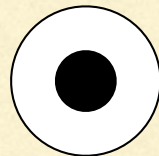
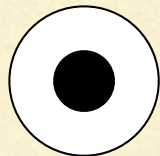
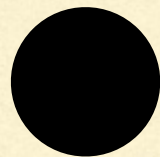
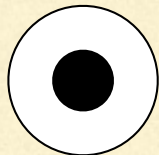
Bayesian inference

Joint probabilities

White,Solid

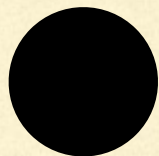


White,Dotted



Black,Dotted

Black,Solid



10 marbles in a bag
Sampling with replacement



$$\Pr(B,S) = 0.4$$



$$\Pr(W,S) = 0.1$$

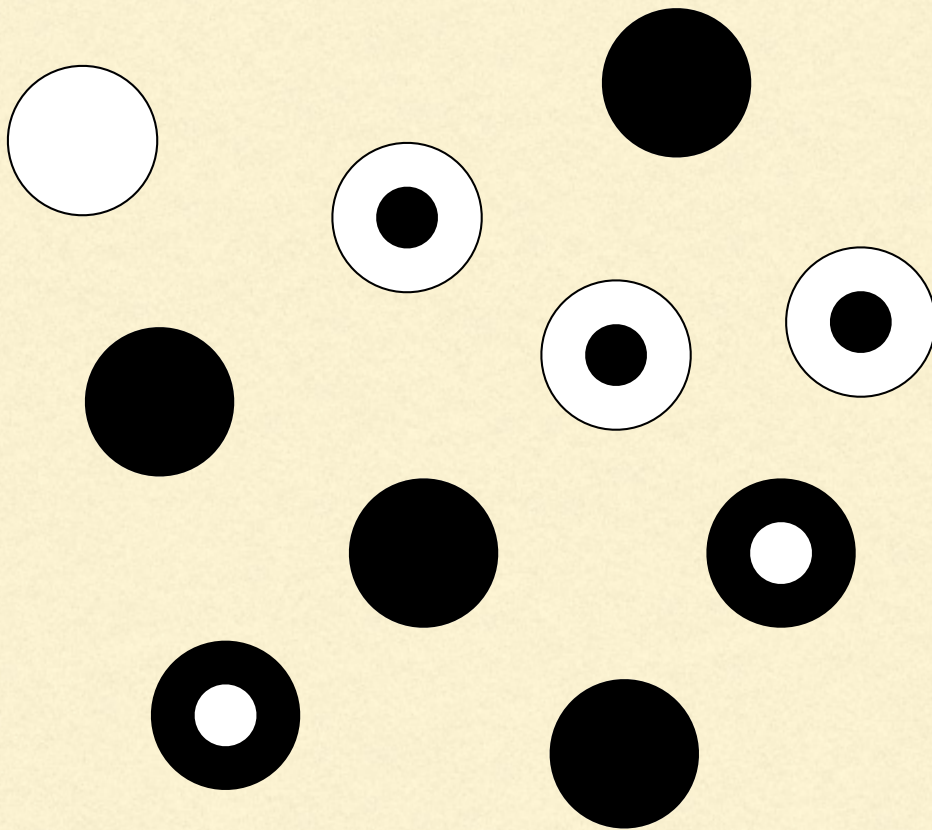


$$\Pr(B,D) = 0.2$$



$$\Pr(W,D) = 0.3$$

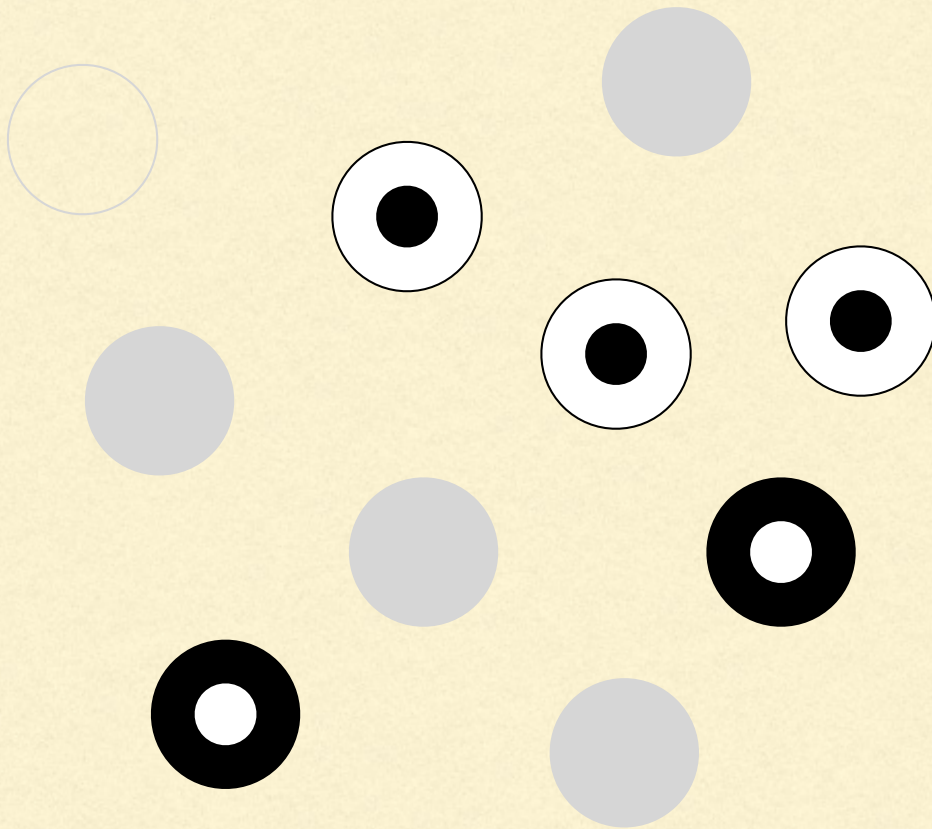
Conditional probabilities



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

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

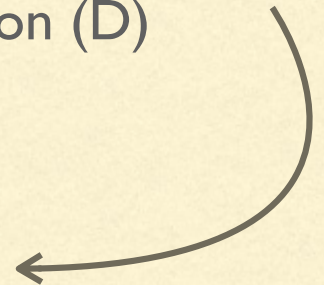
Conditional probabilities



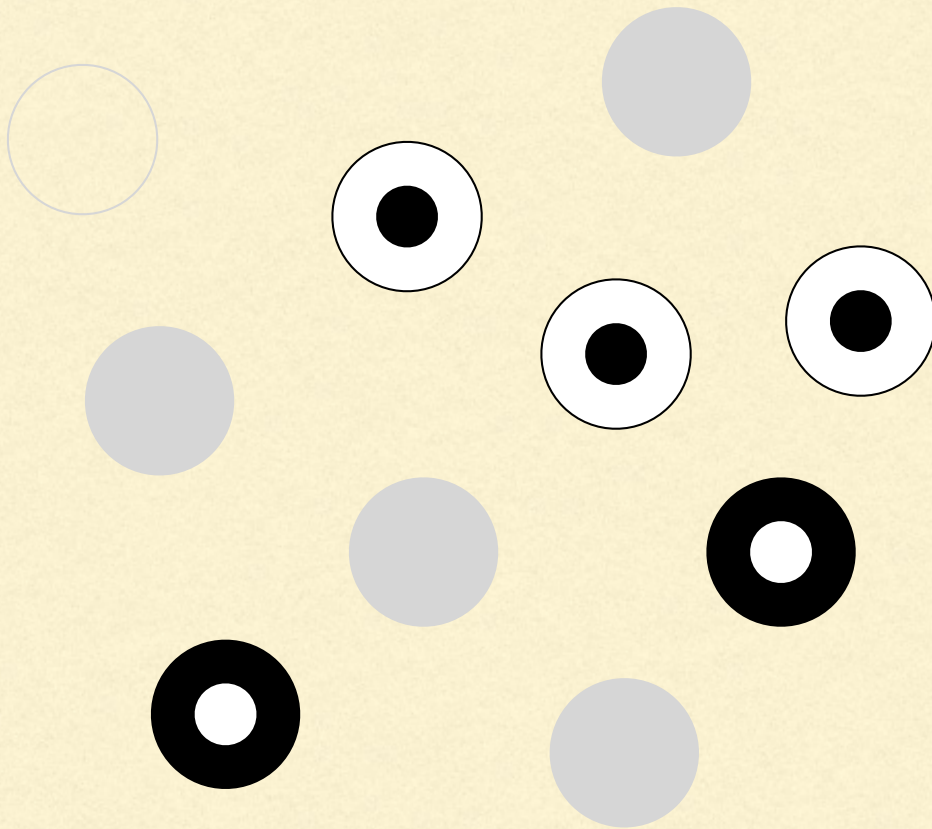
What's the probability that a marble is black given that it is dotted?

5 marbles satisfy the condition (D)

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



Conditional probabilities



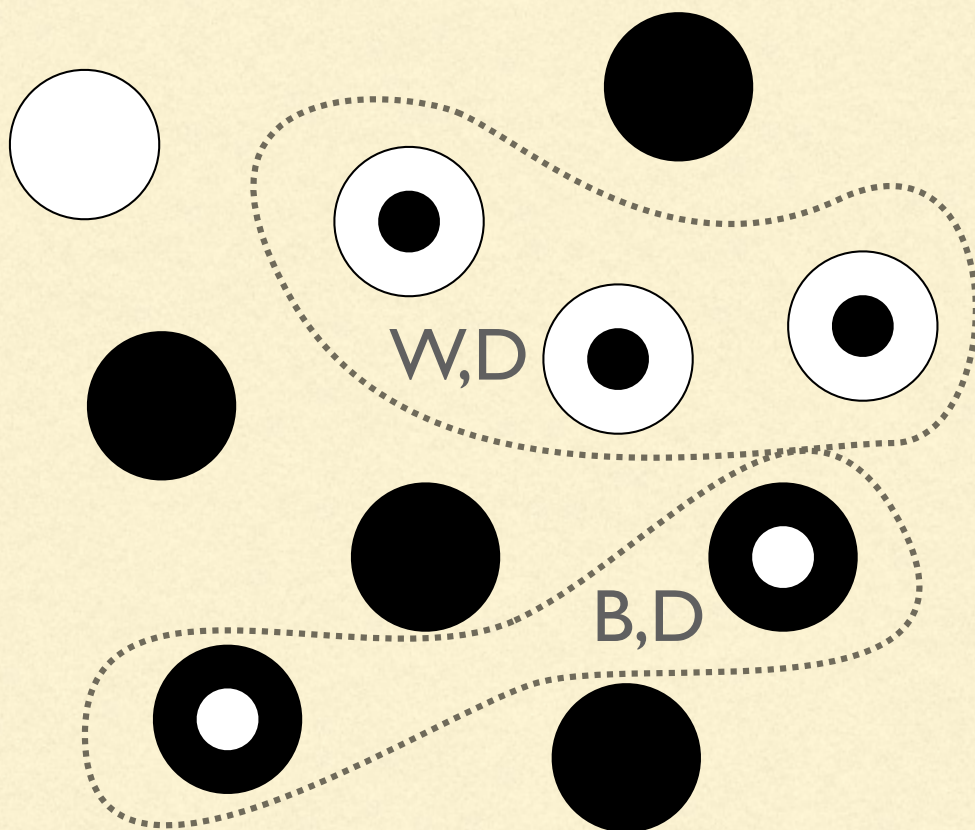
What's the probability that a marble is black given that it is dotted?

5 marbles satisfy the condition (D)

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

2 remaining marbles are black (B)

Marginal probabilities



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

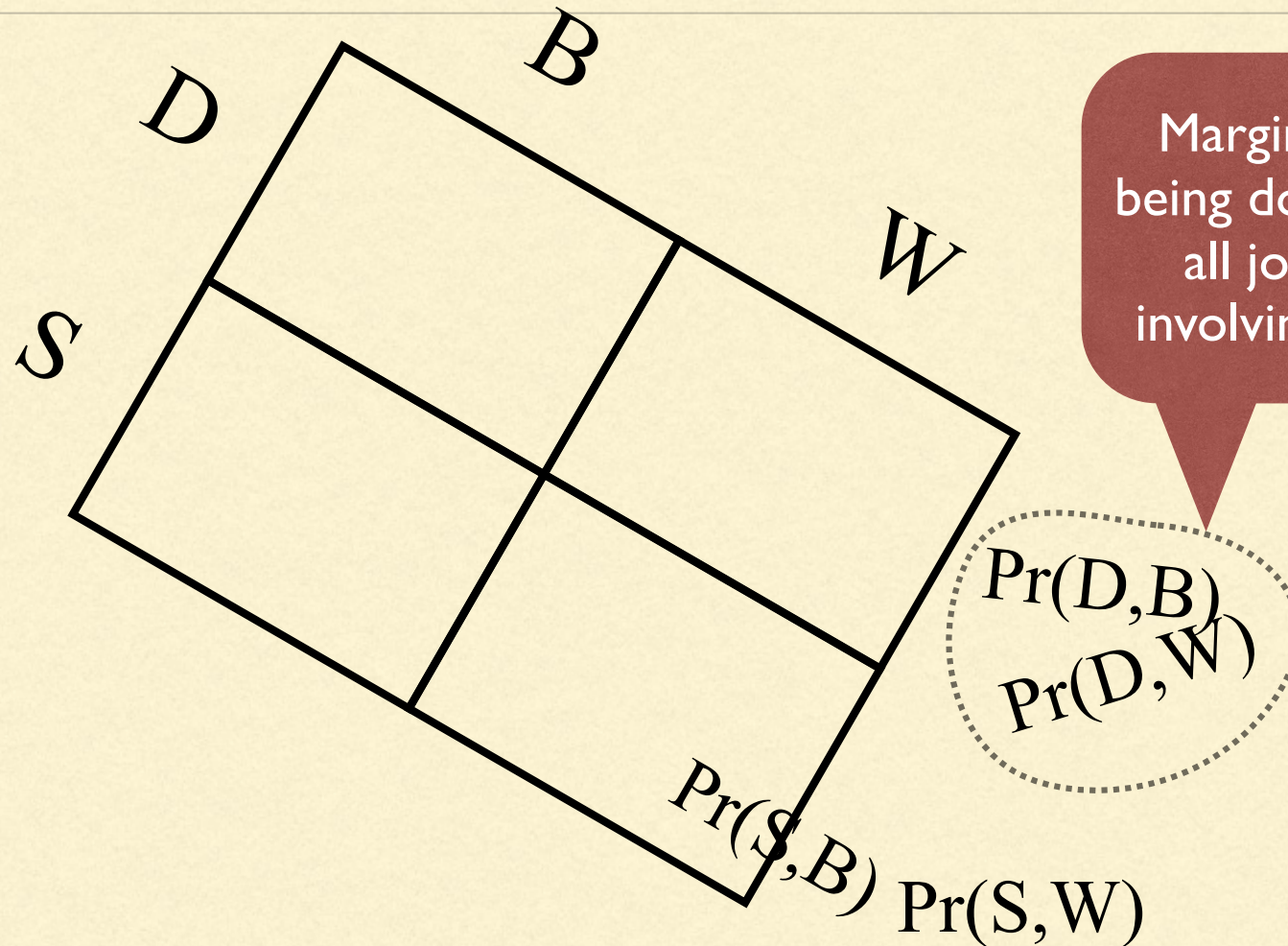
$$\begin{aligned}\Pr(\mathbf{D}) &= \Pr(\mathbf{B}, \mathbf{D}) + \Pr(\mathbf{W}, \mathbf{D}) \\ &= 0.2 + 0.3 \\ &= 0.5\end{aligned}$$

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



Marginal probabilities

	B	W
D		
S		

$\Pr(D)$ = marginal probability
of being dotted

$$\Pr(D,B) + \Pr(D,W)$$

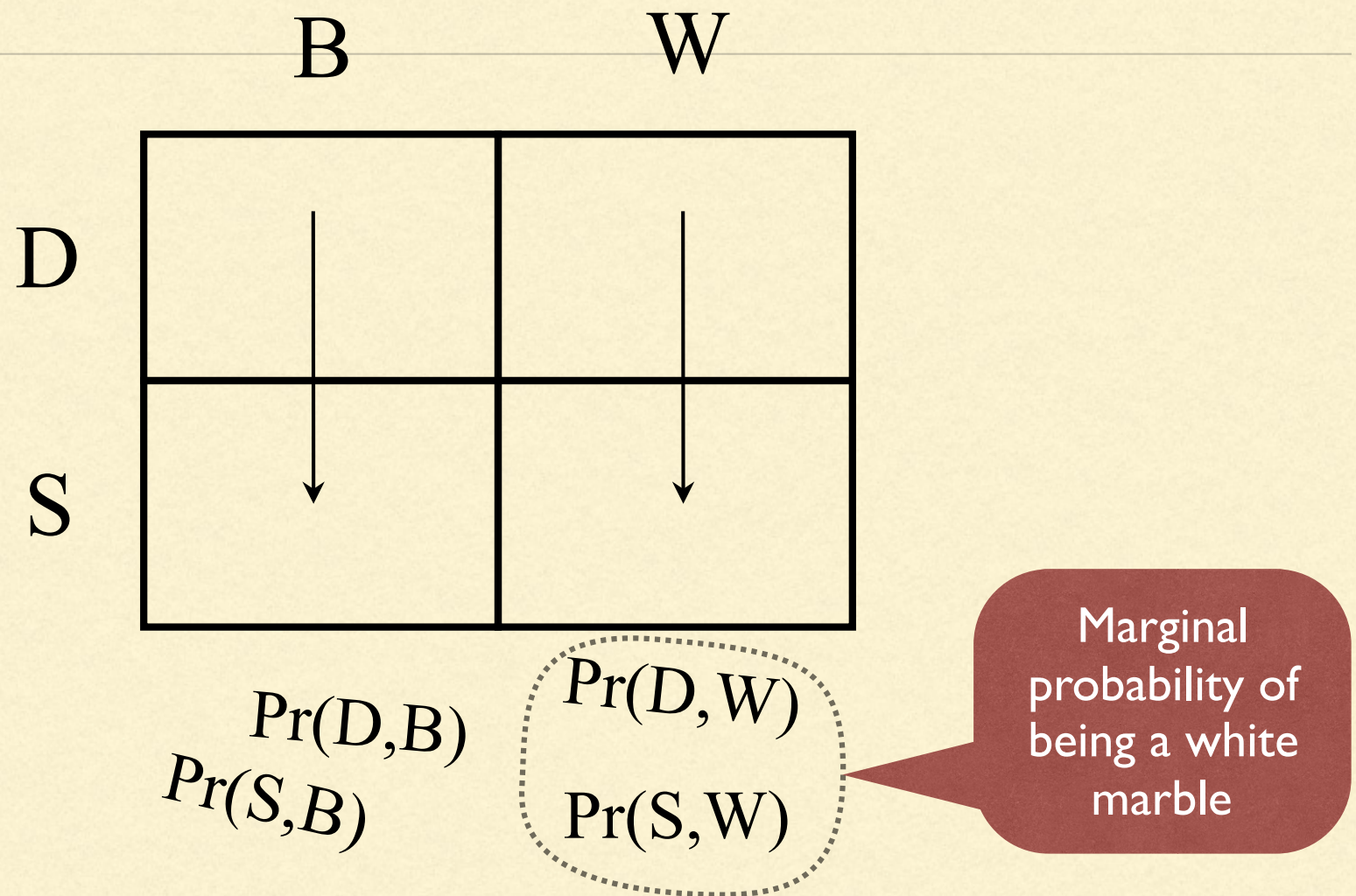
$$\Pr(S,B) + \Pr(S,W)$$

$\Pr(S)$ = marginal probability
of being solid

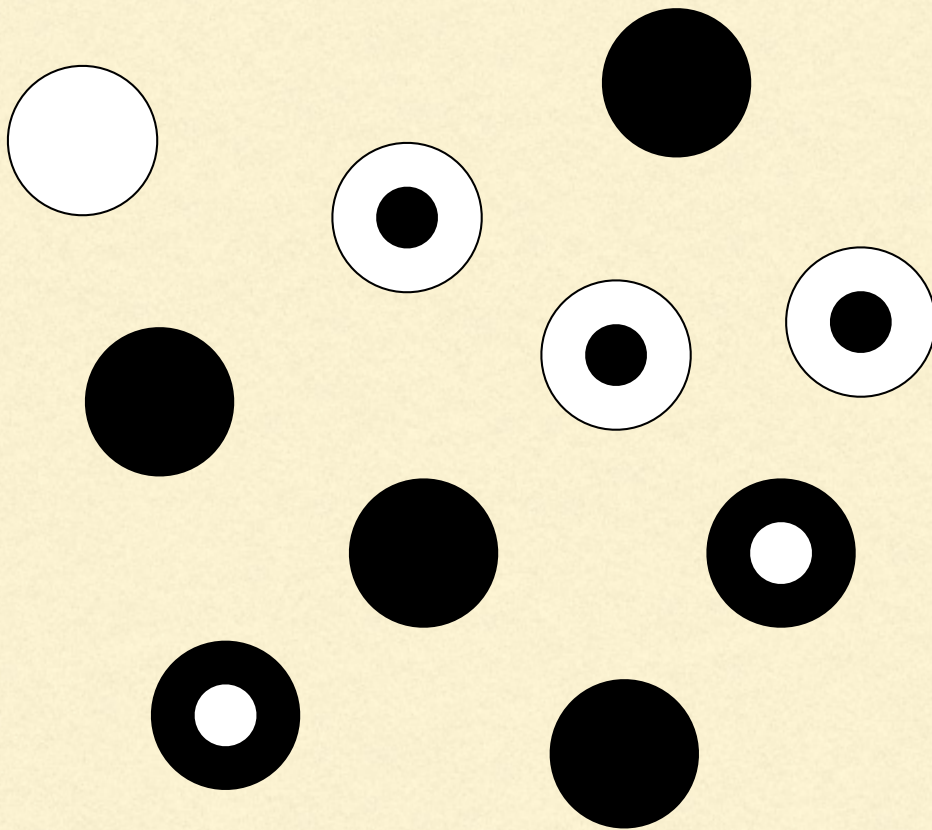
Joint probabilities

	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

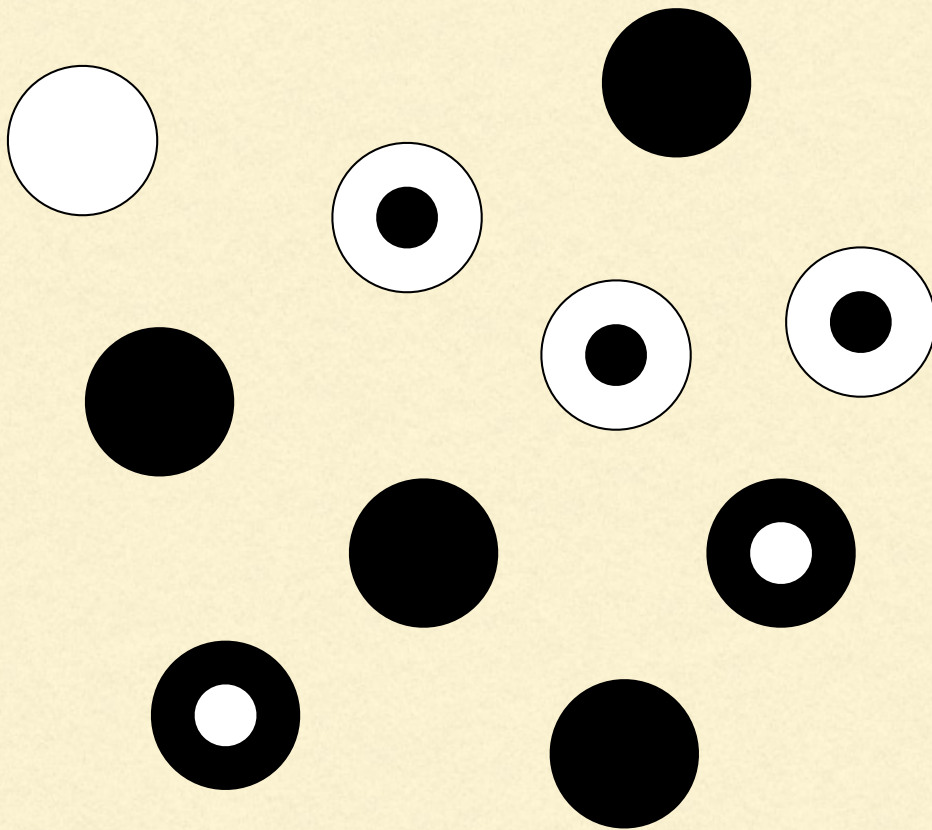
$$\Pr(B,D) \begin{cases} \Pr(B|D) \Pr(D) \\ \Pr(D|B) \Pr(B) \end{cases}$$

Either B or D
can be the
condition

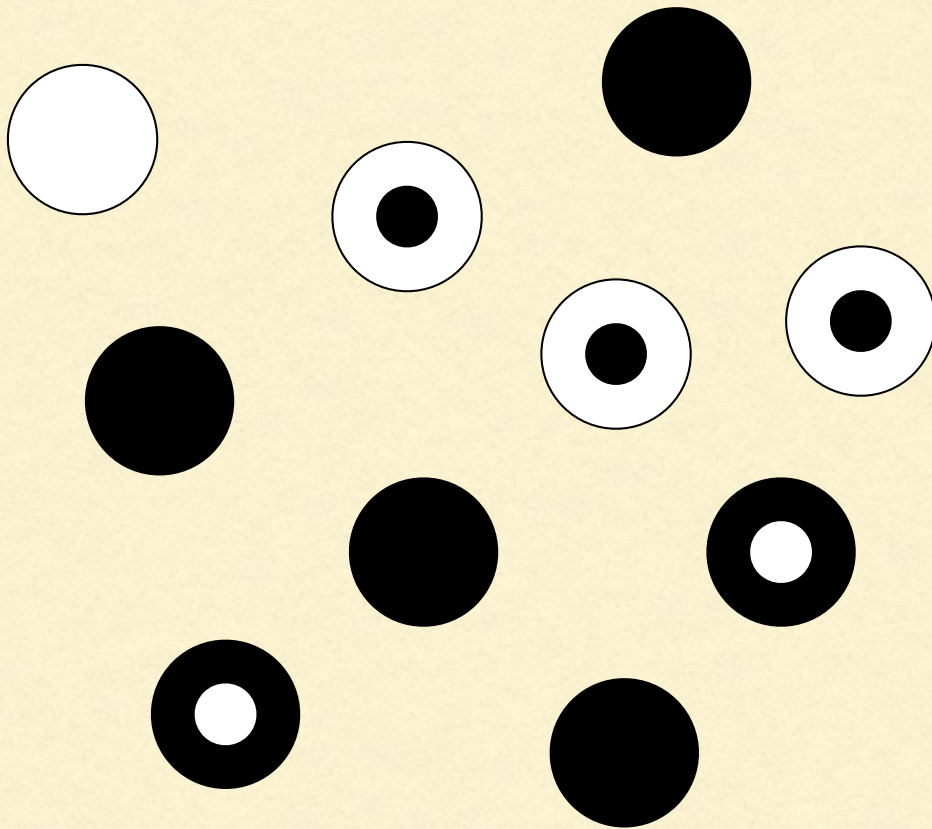
Bayes' rule

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

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



Bayes' rule



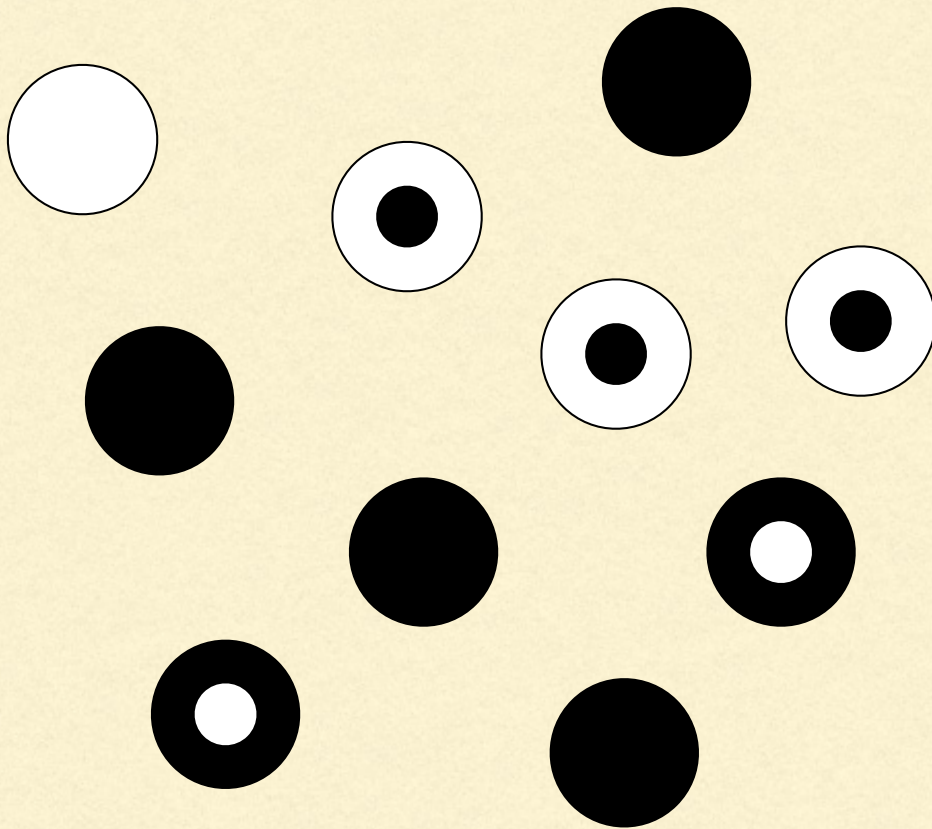
Equate the two ways of writing $\Pr(B,D)$

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

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



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

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

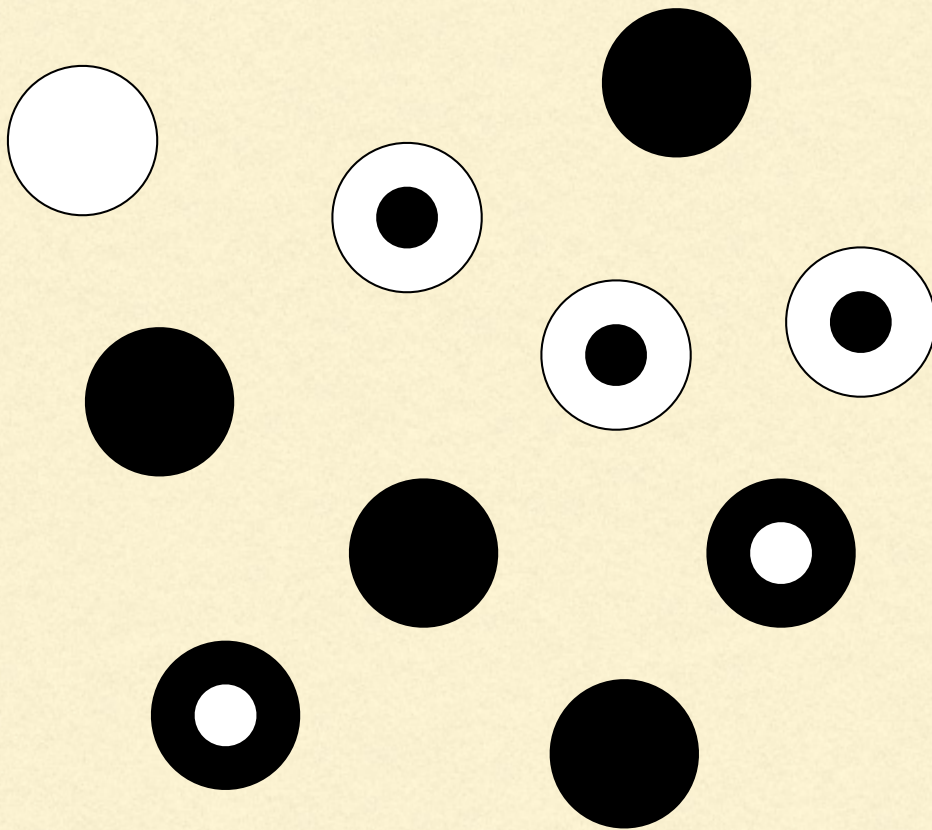
Divide both sides by $\Pr(D)$

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

Bayes' rule

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

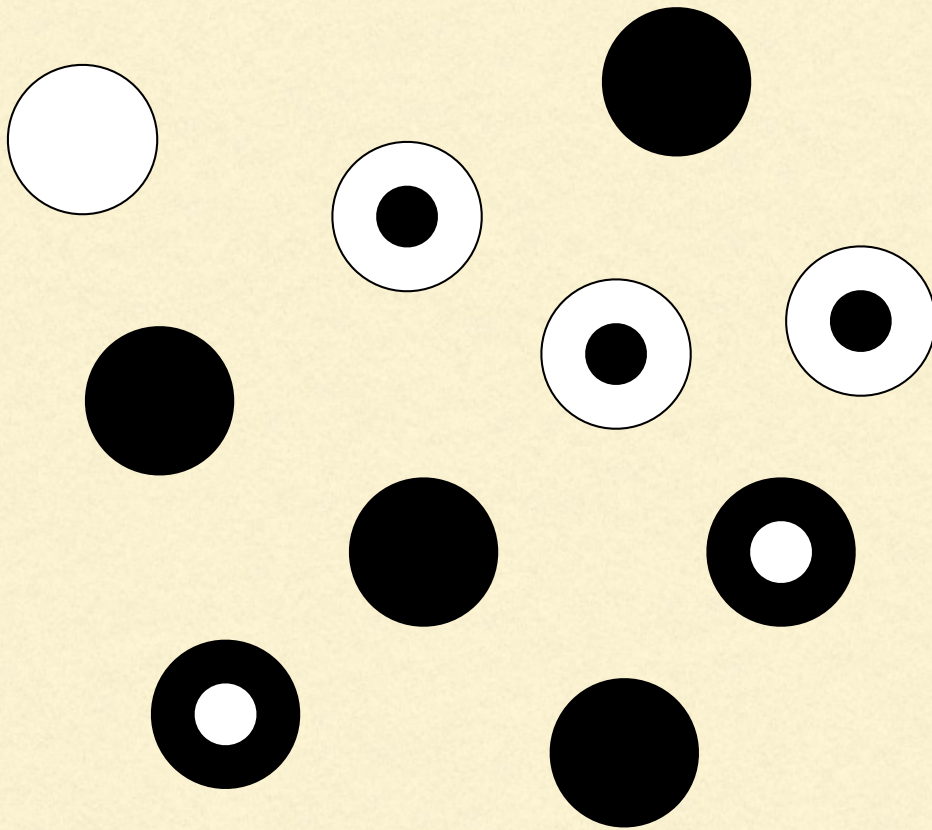
Bayes' rule



Bayes' rule

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

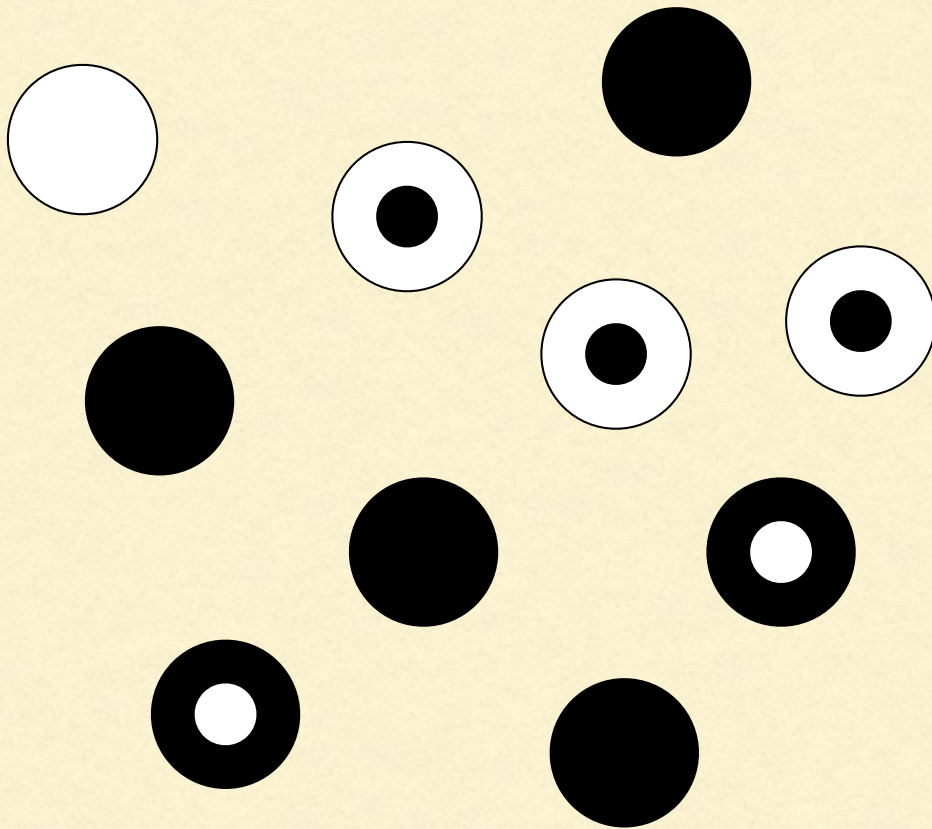
Bayes' rule



Bayes' rule

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

Bayes' rule

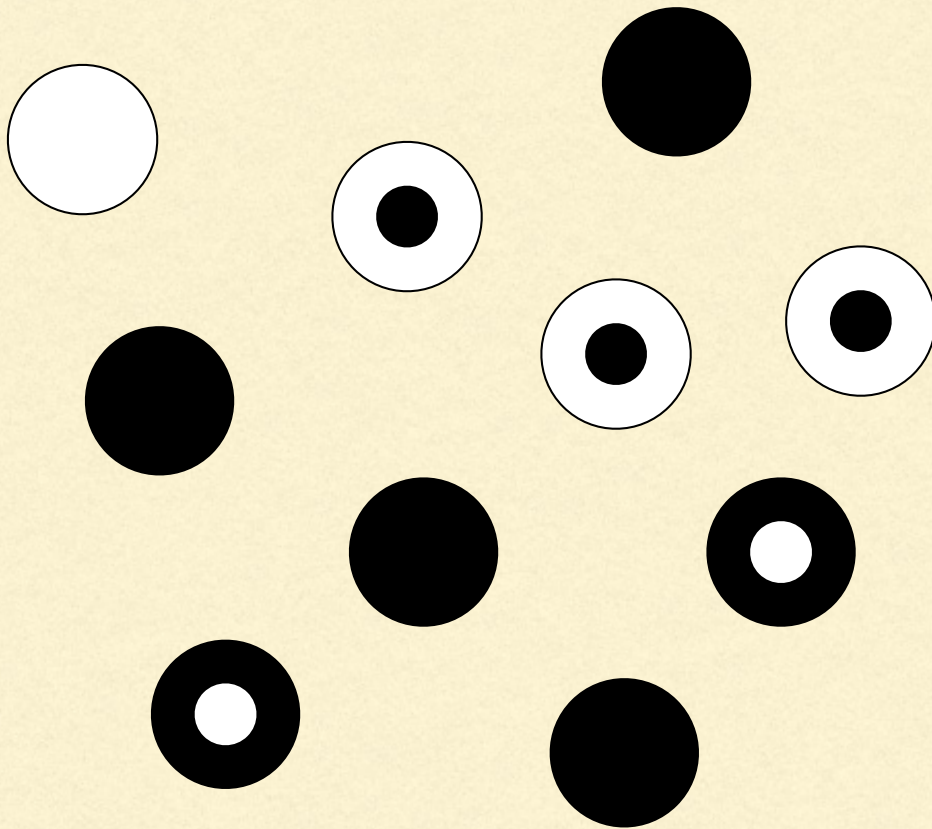


Bayes' rule

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

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

Bayes' rule



Bayes' rule

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

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

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



Bayes' rule (variations)

$$\begin{aligned}\Pr(B|D) &= \frac{\Pr(B) \Pr(D|B)}{\Pr(D)} \\ &= \frac{\Pr(D, B)}{\Pr(D, B) + \Pr(D, W)}\end{aligned}$$

Bayes' rule (variations)

$$\begin{aligned}\Pr(B|D) &= \frac{\Pr(B) \Pr(D|B)}{\Pr(D)} \\ &= \frac{\Pr(D, B)}{\Pr(D, B) + \Pr(D, W)}\end{aligned}$$

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

Normalizing constant

It is easy to see that $\Pr(D)$ serves as a *normalization constant*, ensuring that $\Pr(B|D) + \Pr(W|D) = 1.0$

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

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

$$\Pr(B|D) + \Pr(W|D) = \frac{\cancel{\Pr(D, B)} + \cancel{\Pr(D, W)}}{\cancel{\Pr(D, B)} + \cancel{\Pr(D, W)}} = 1$$

Bayes' rule (variations)

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

Bayes' rule (variations)

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

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

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 θ

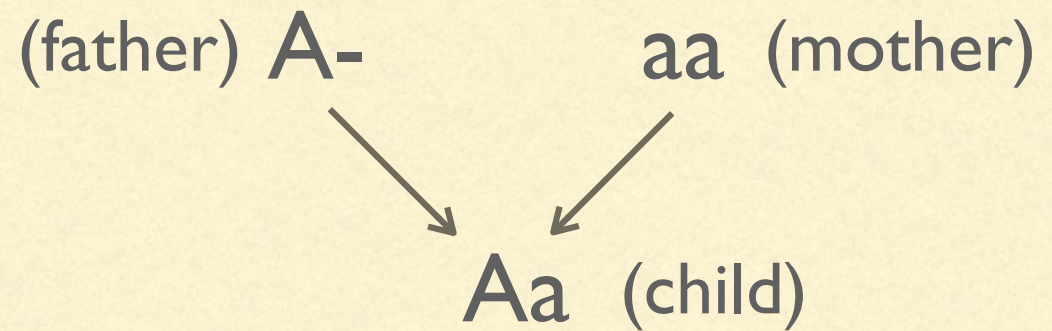
Prior probability of hypothesis θ

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

Posterior probability of hypothesis θ

Marginal probability of the data (marginalizing over hypotheses)

Paternity example



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

θ_1

θ_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*

The diagram shows the equation for Bayes' rule in the continuous case. The numerator consists of two terms: $p(D | \theta)$ (highlighted in a light blue box) and $p(\theta)$ (highlighted in a light red box). The denominator is the integral $\int p(D | \theta) p(\theta) d\theta$ (highlighted in a light green box). The entire equation is enclosed in a light purple box. Arrows point from the labels 'Likelihood' and 'Prior probability density' to the blue and red boxes respectively. An arrow points from the label 'Posterior probability density' to the purple box. Another arrow points from the label 'Marginal probability of the data' to the green box.

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

Posterior probability *density*

Marginal probability of the data

If you had to guess...

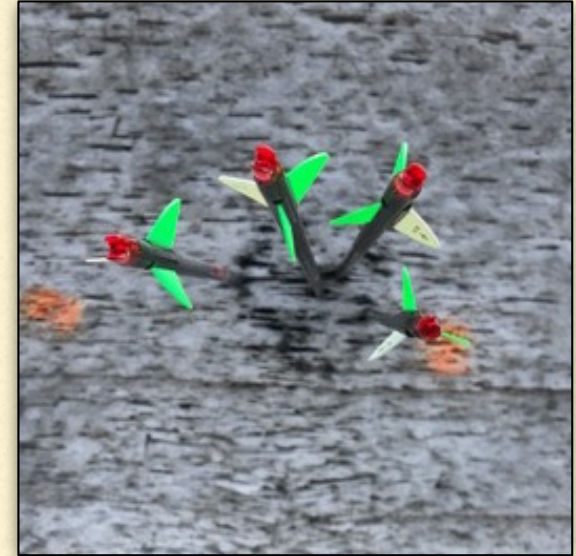
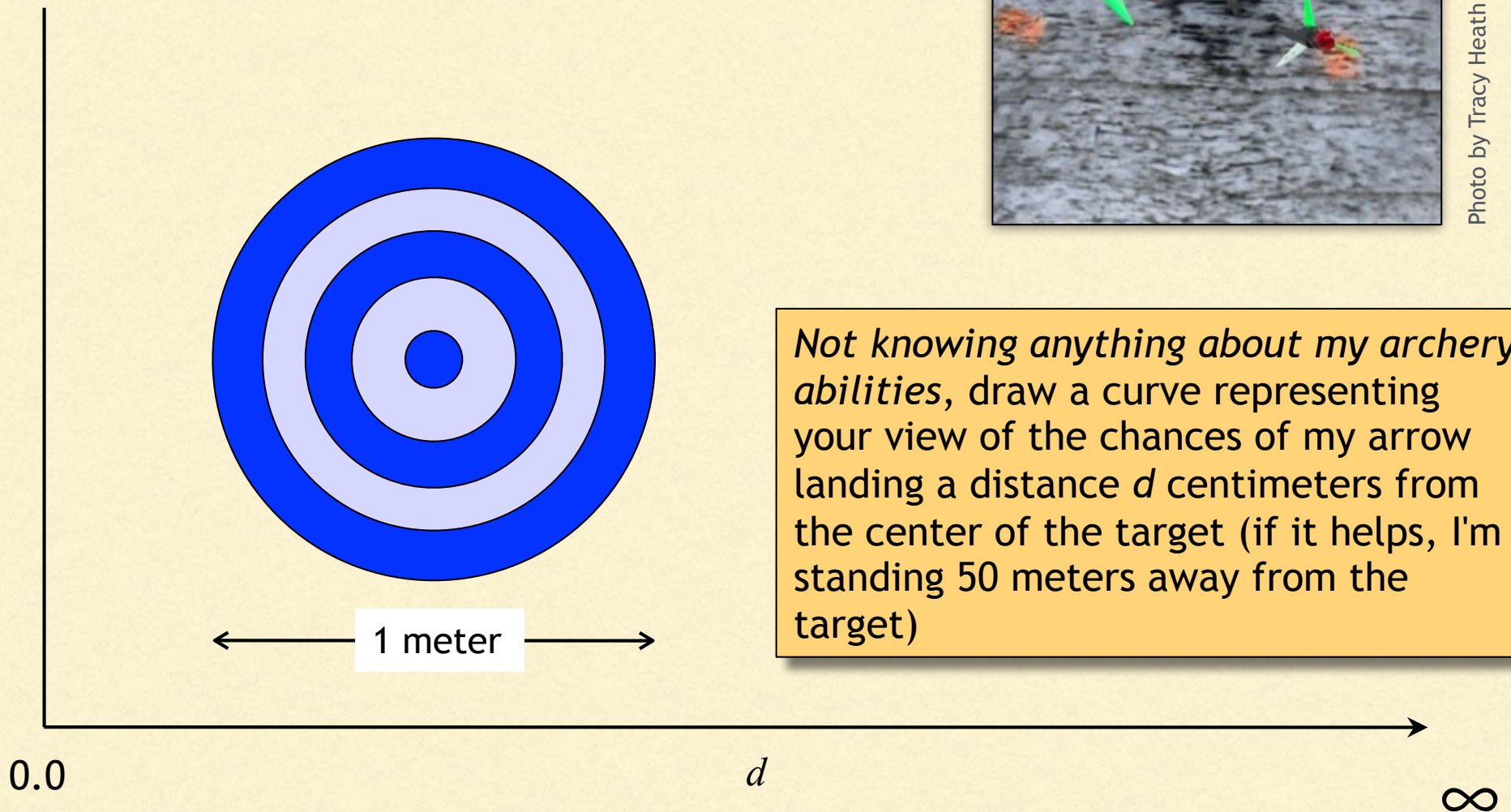
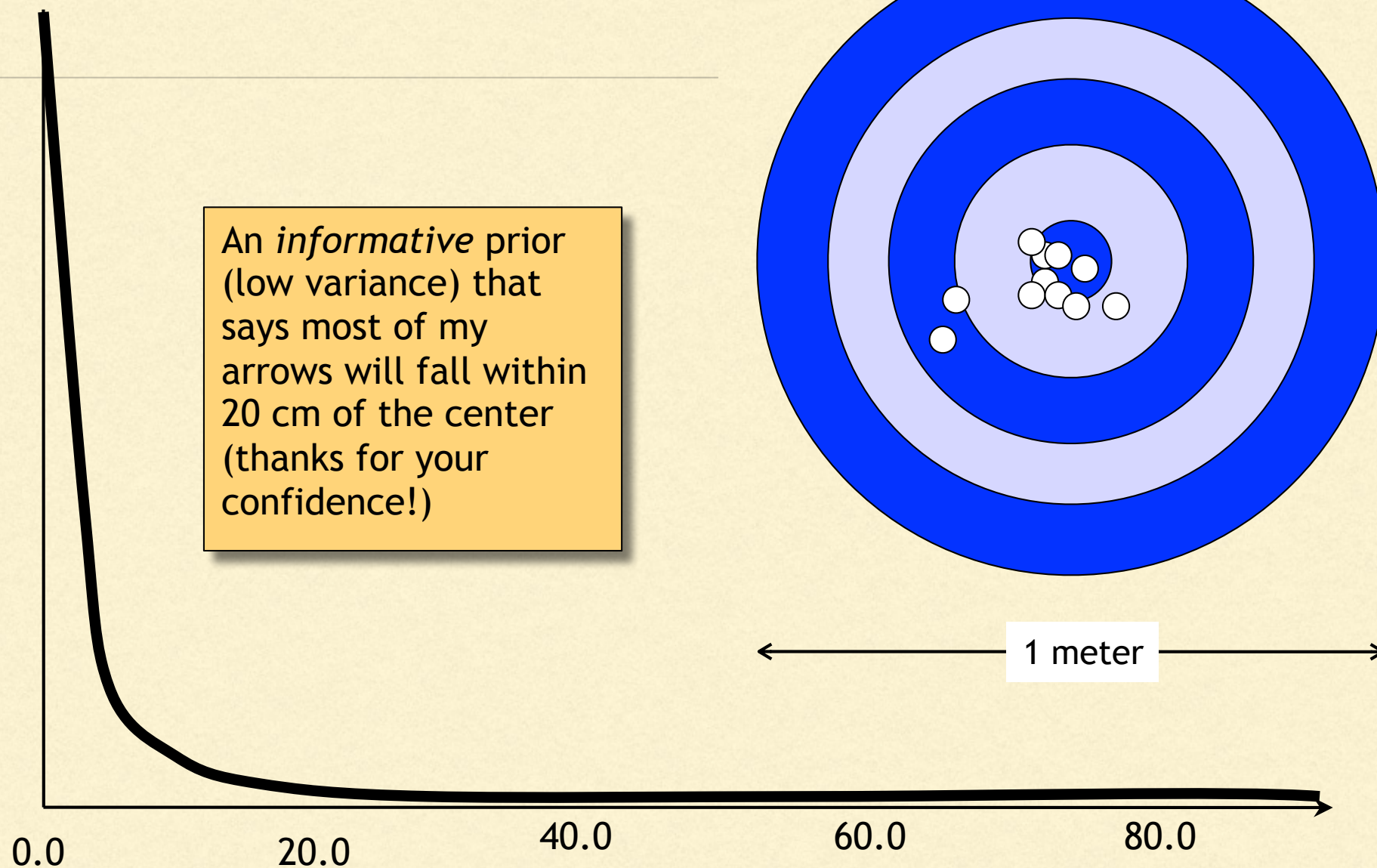


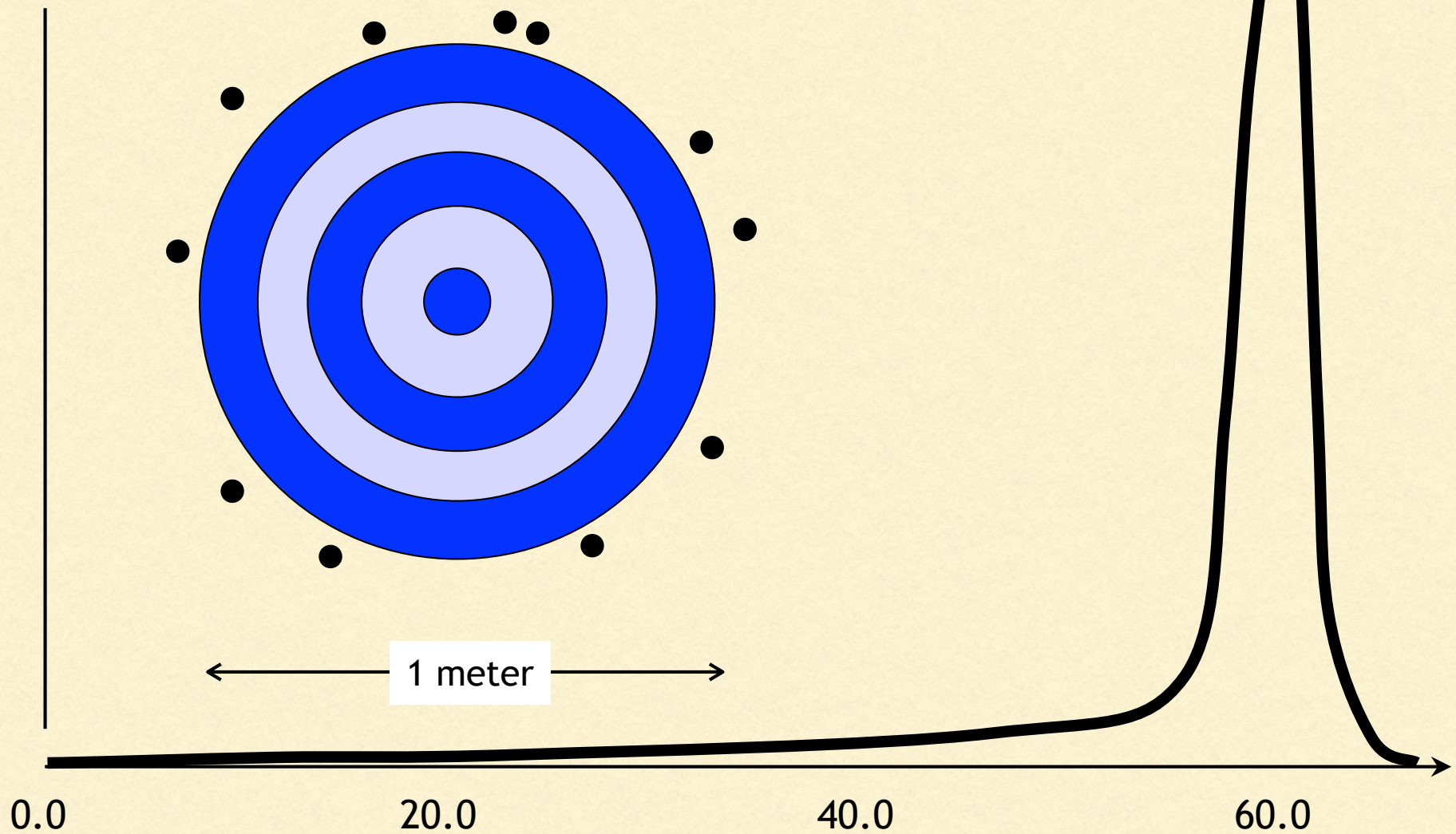
Photo by Tracy Heath



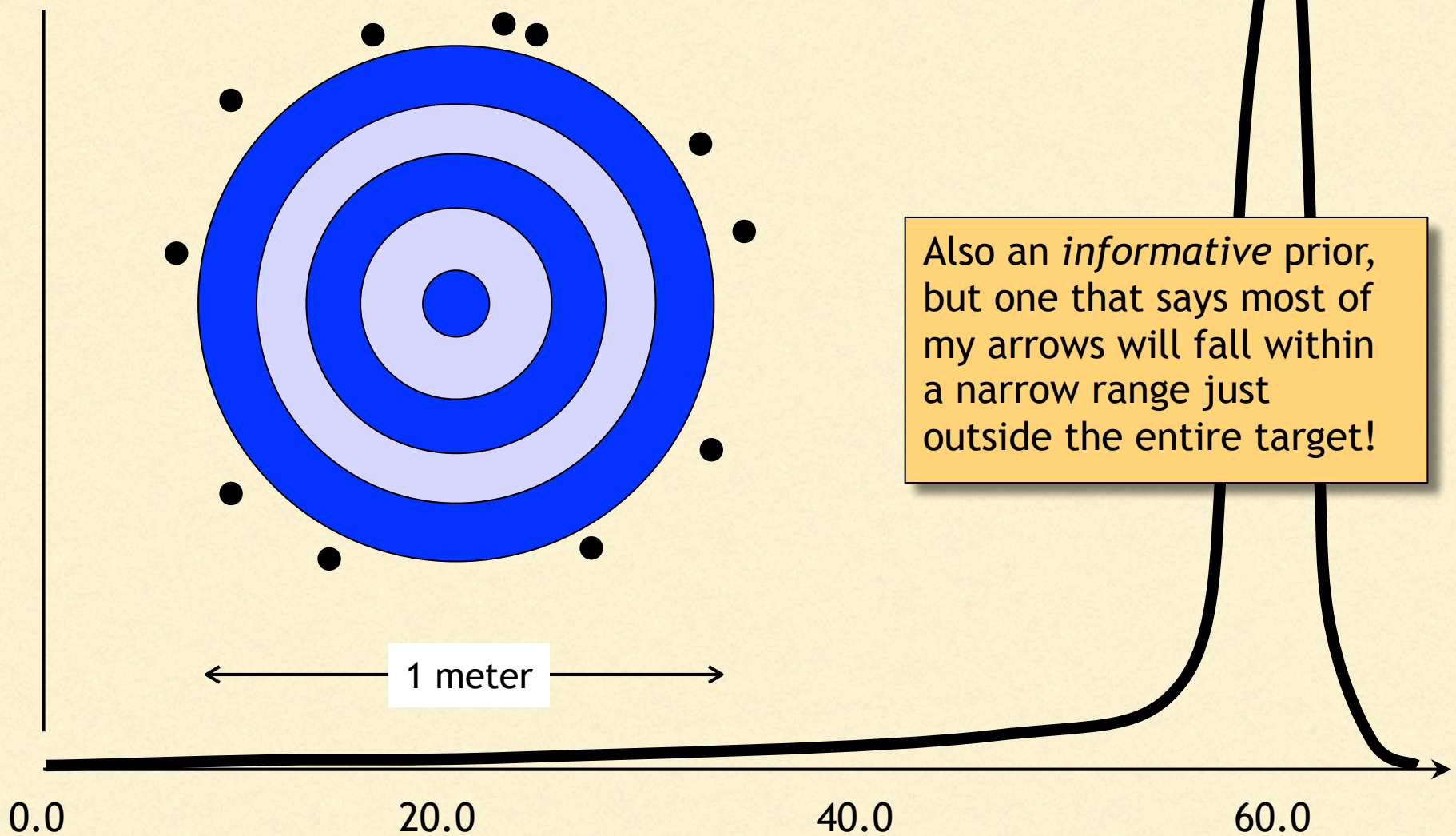
Case 1: assume I have talent



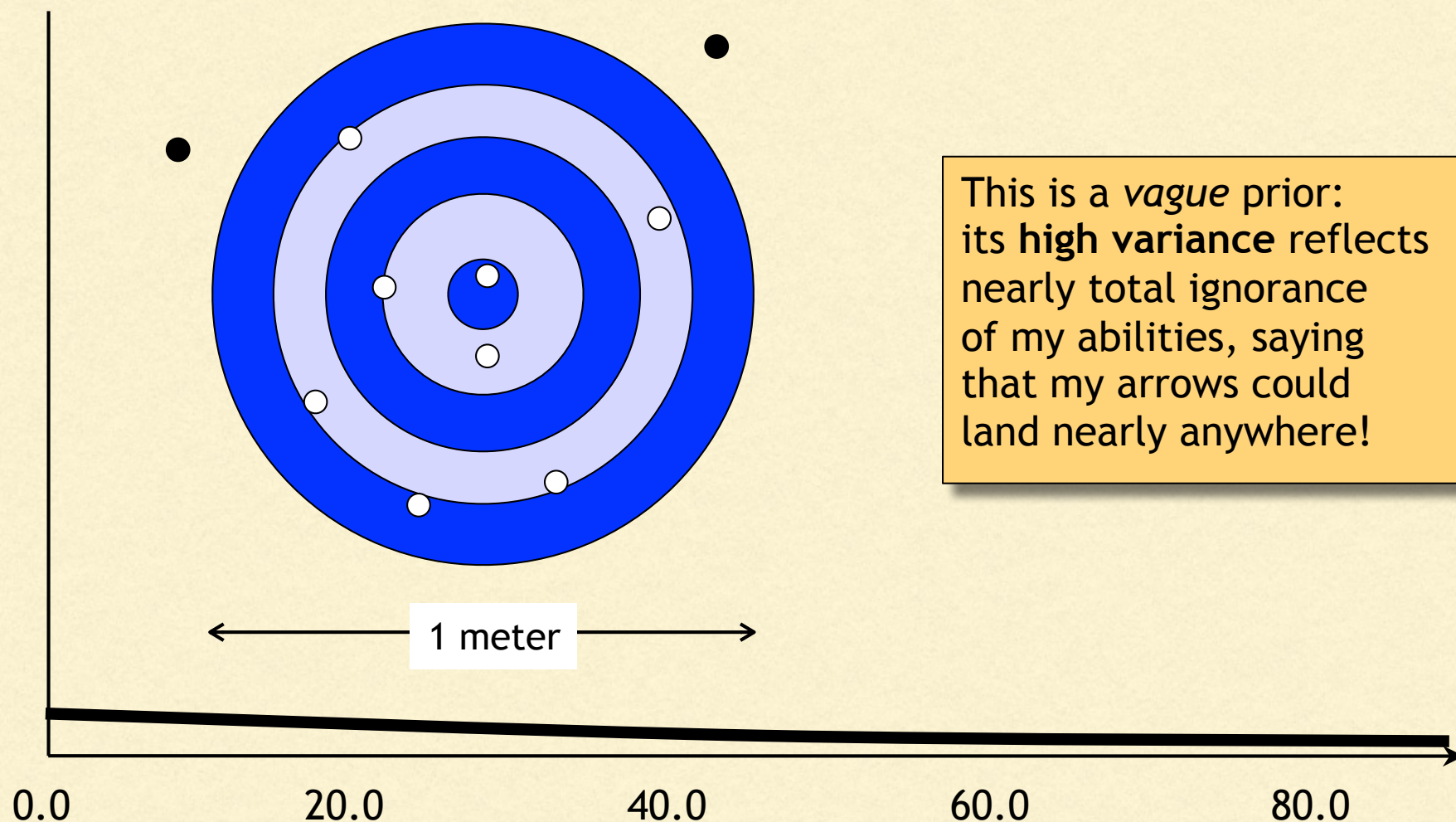
Case 2: assume I have a talent for missing the target!



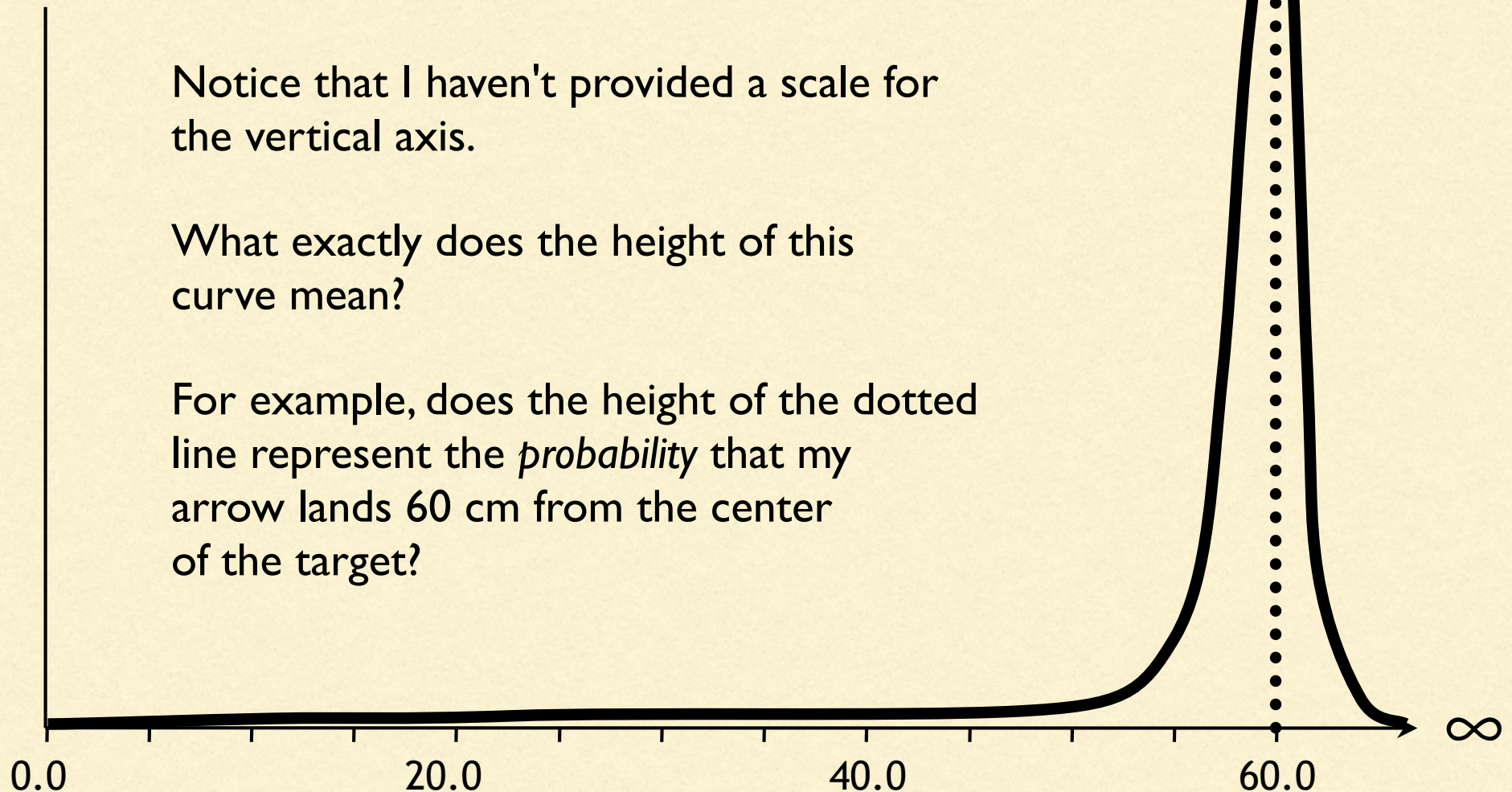
Case 2: assume I have a talent for missing the target!



Case 3: assume I have no talent



A matter of scale

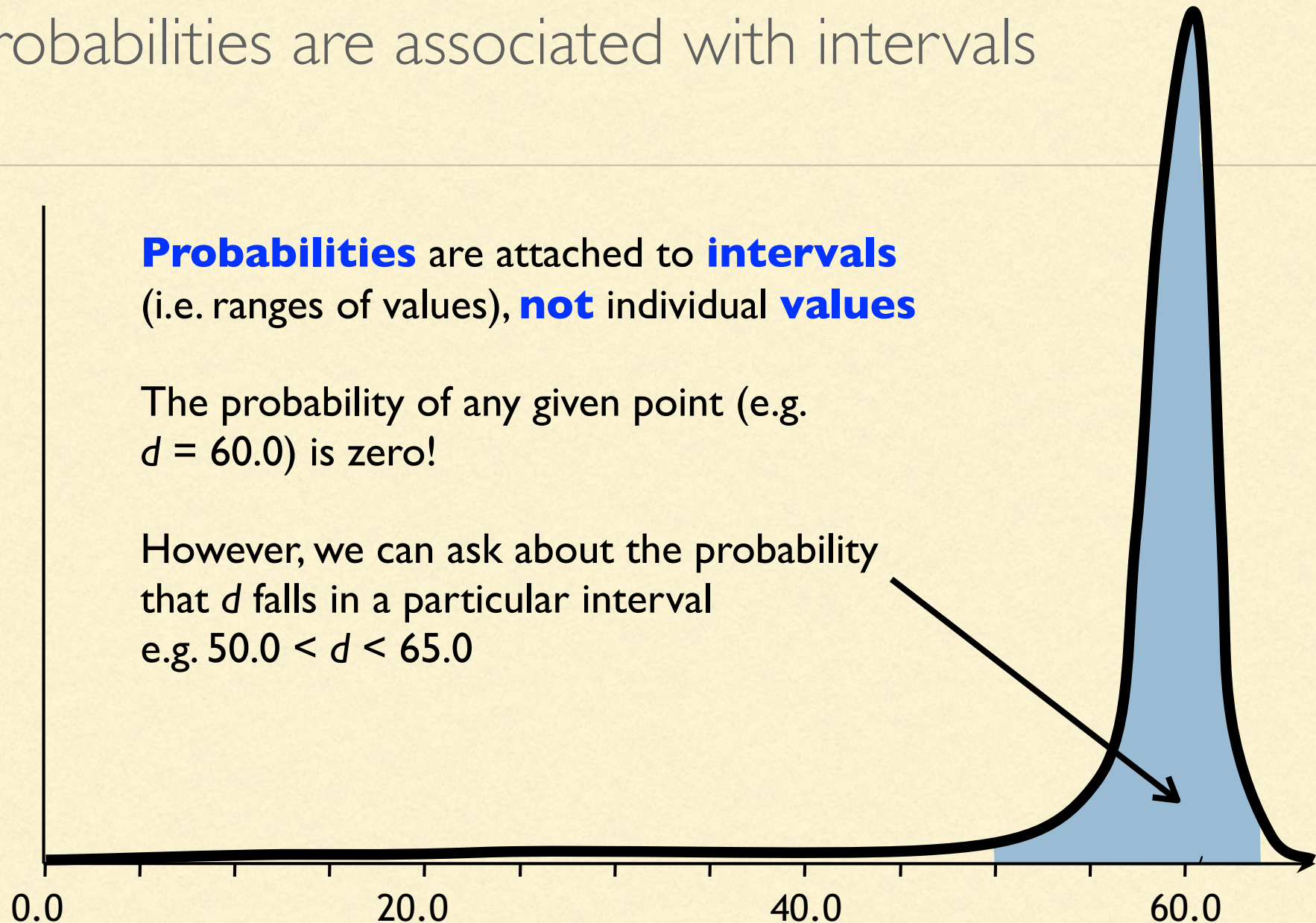


Probabilities are associated with intervals

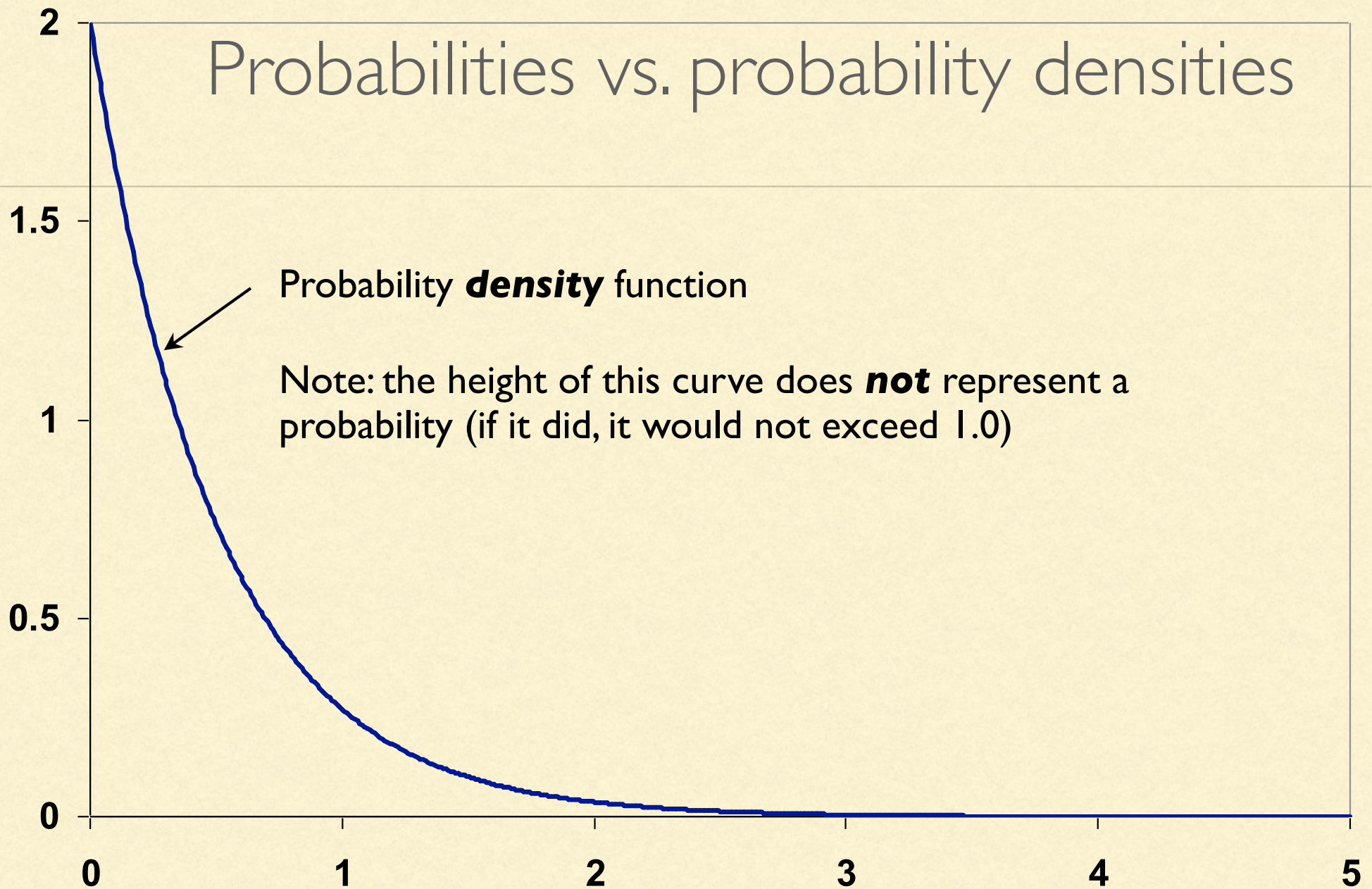
Probabilities are attached to **intervals**
(i.e. ranges of values), **not** individual **values**

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. $50.0 < d < 65.0$



Probabilities vs. probability densities

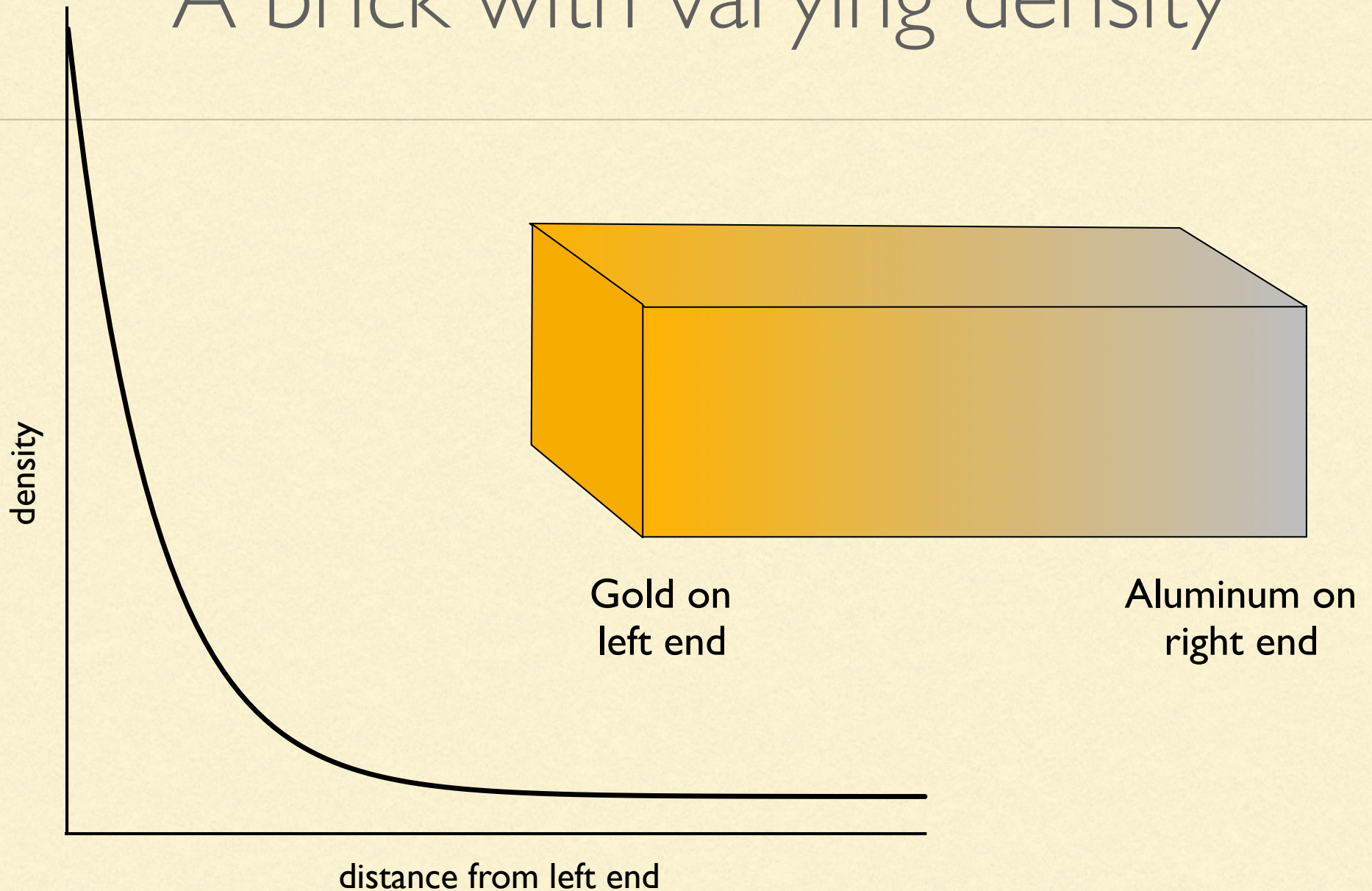


Densities of various substances

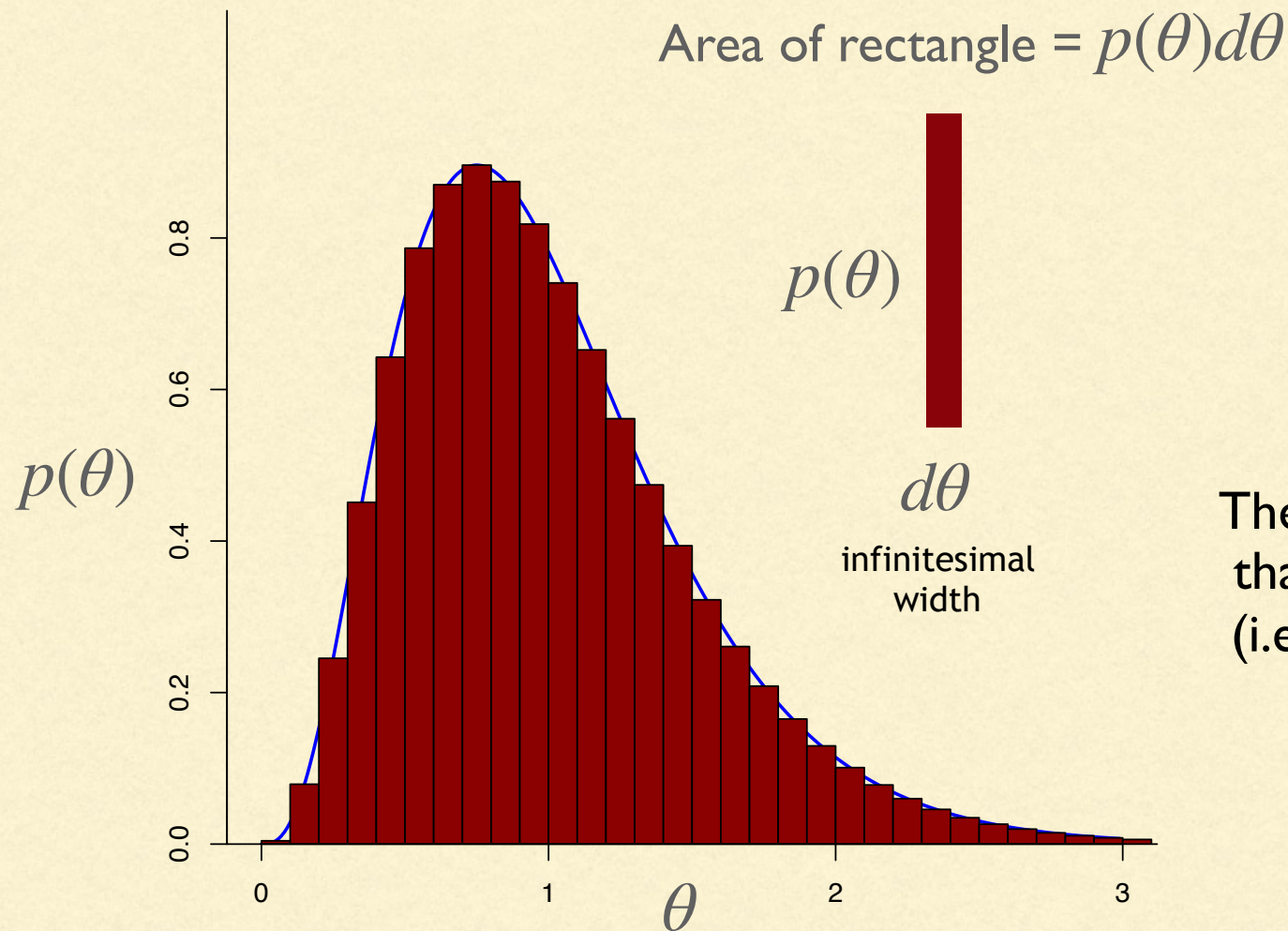
Substance	Density (g/cm ³)
Cork	0.24
Aluminum	2.7
Gold	19.3

Density does not equal mass
 $\text{mass} = \text{density} \times \text{volume}$

A brick with varying density



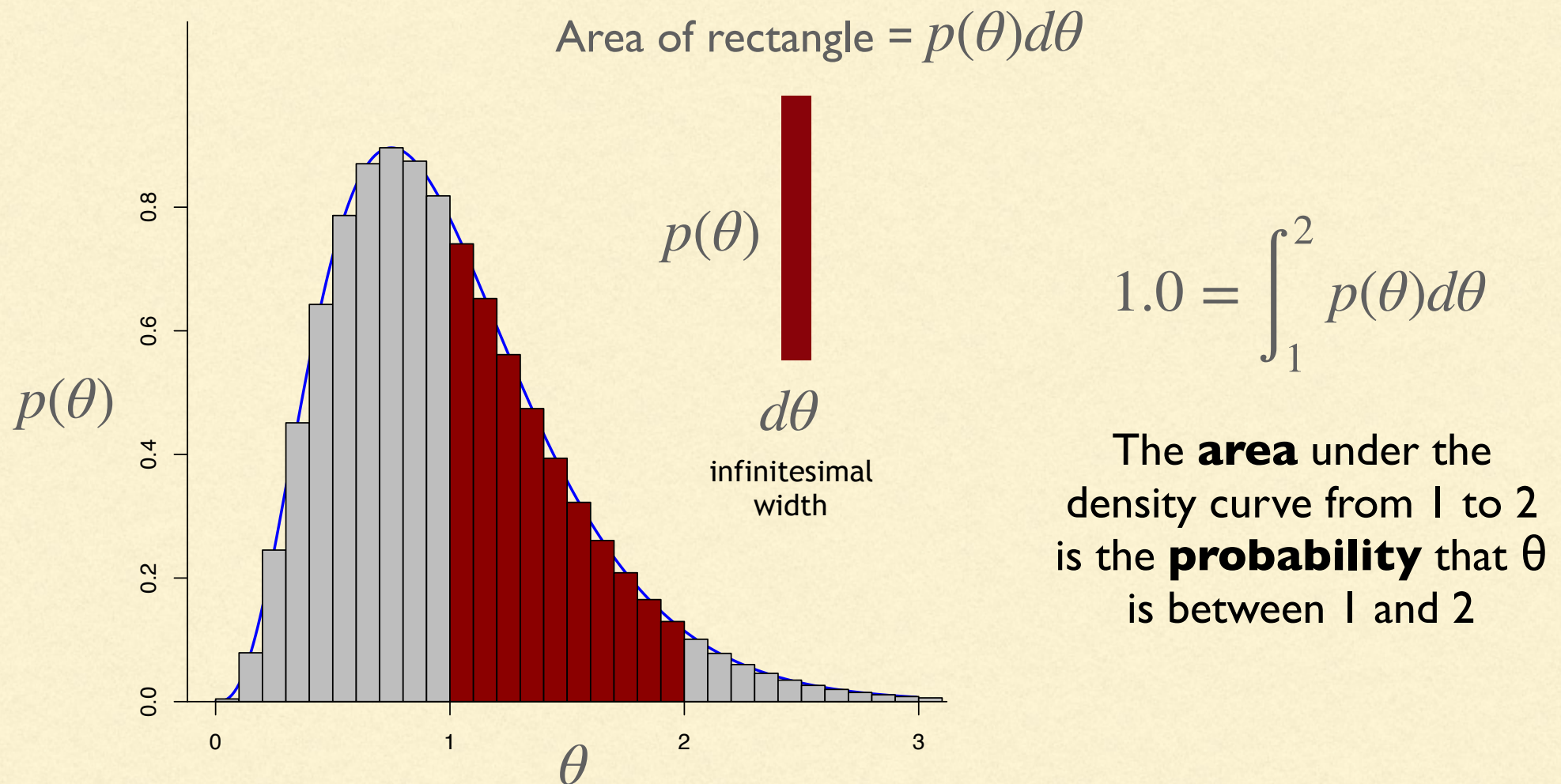
Integrating a density yields a probability



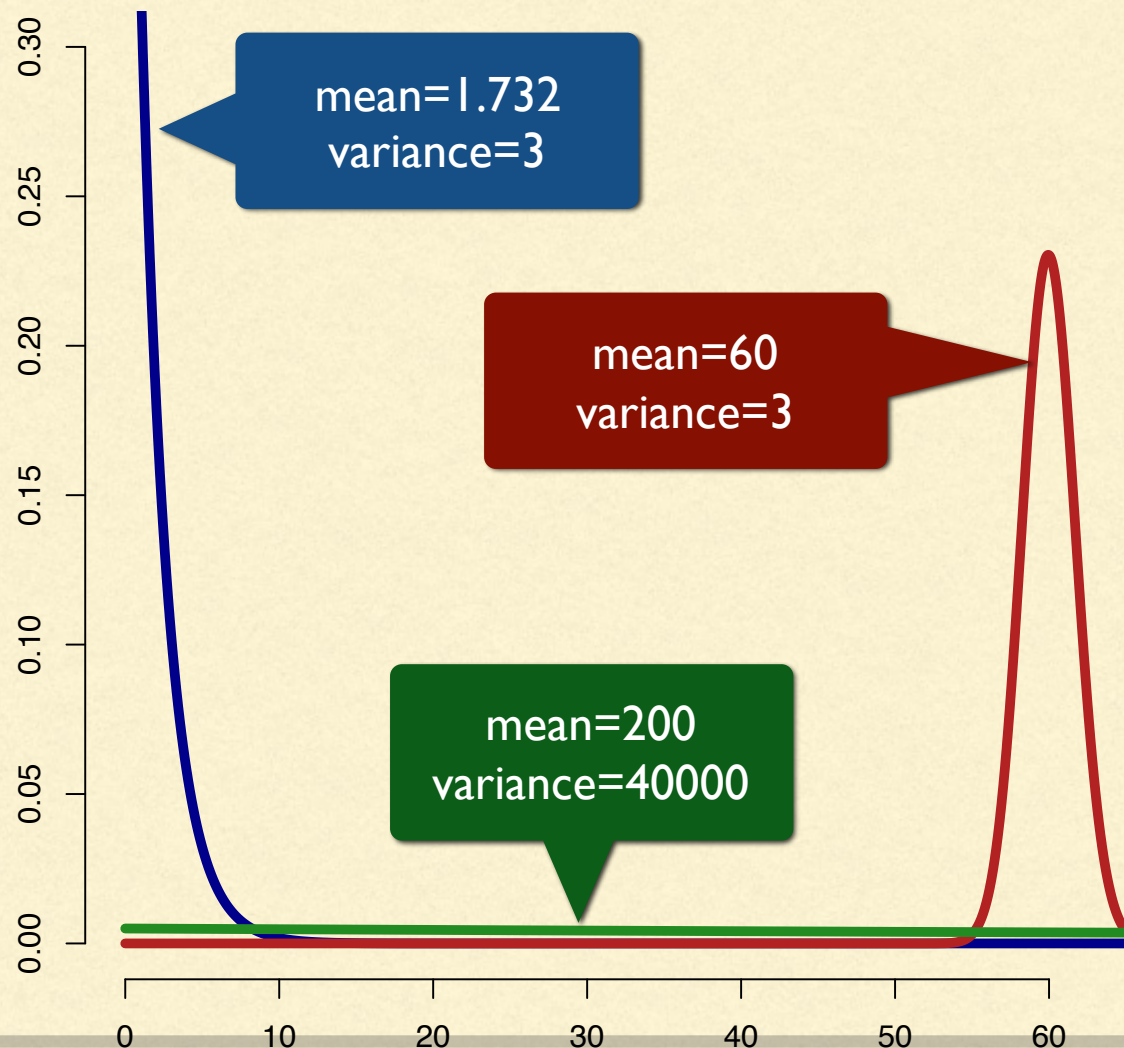
$$1.0 = \int p(\theta)d\theta$$

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

Integrating a density yields a probability



Archery priors revisited



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**.

Archery Prior Applet

<https://phylogeny.uconn.edu/archery-priors/>

Usually there are many parameters...

A 2-parameter example

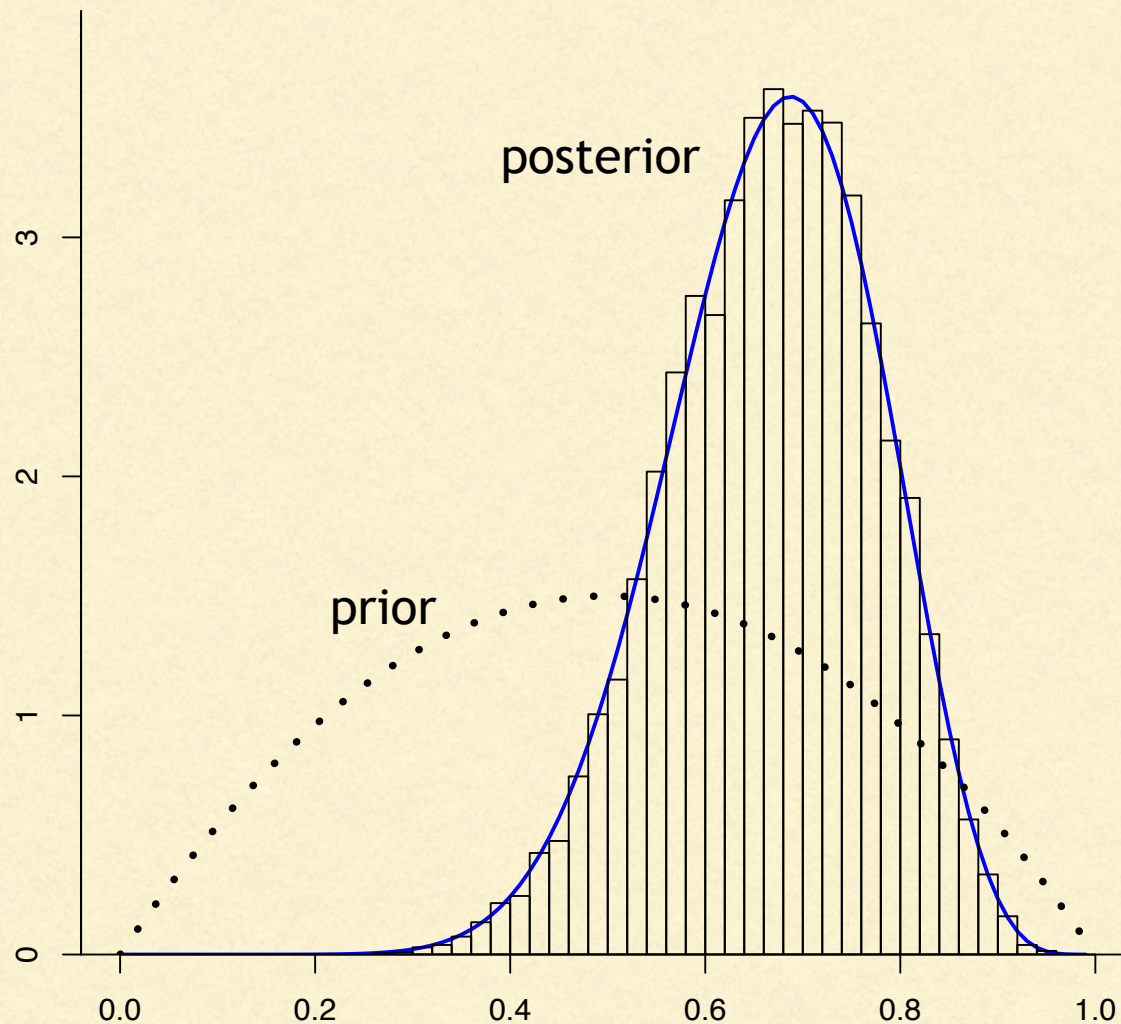
$$p(\theta, \phi | D) = \frac{\overbrace{p(D | \theta, \phi) p(\theta) p(\phi)}^{\text{Likelihood} \quad \text{Prior density}}}{\underbrace{\int_{\theta} \int_{\phi} p(D | \theta, \phi) p(\theta) p(\phi) d\phi d\theta}_{\text{Marginal probability of data}}}$$

↑
Posterior
probability
density

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...

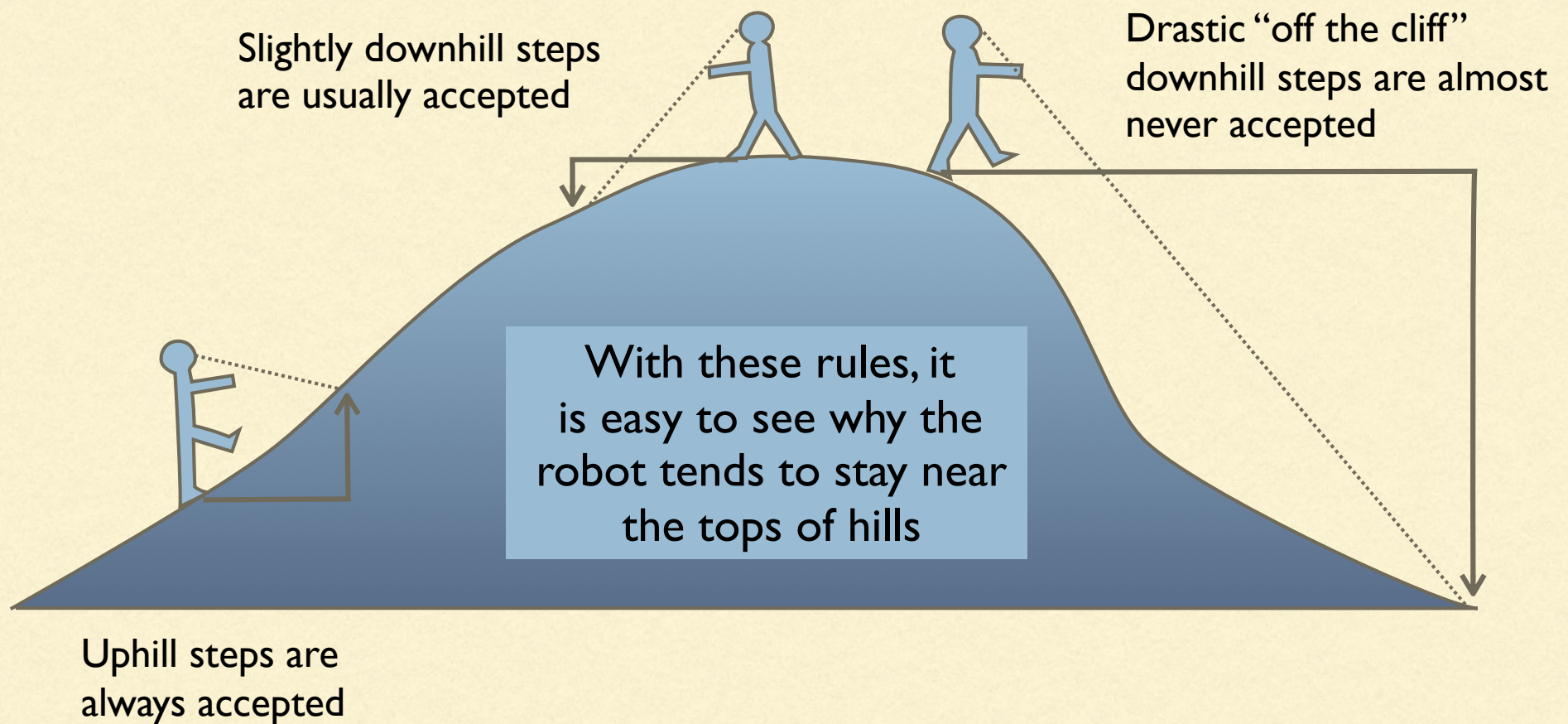
Markov chain Monte Carlo (MCMC)

Markov chain Monte Carlo (MCMC)

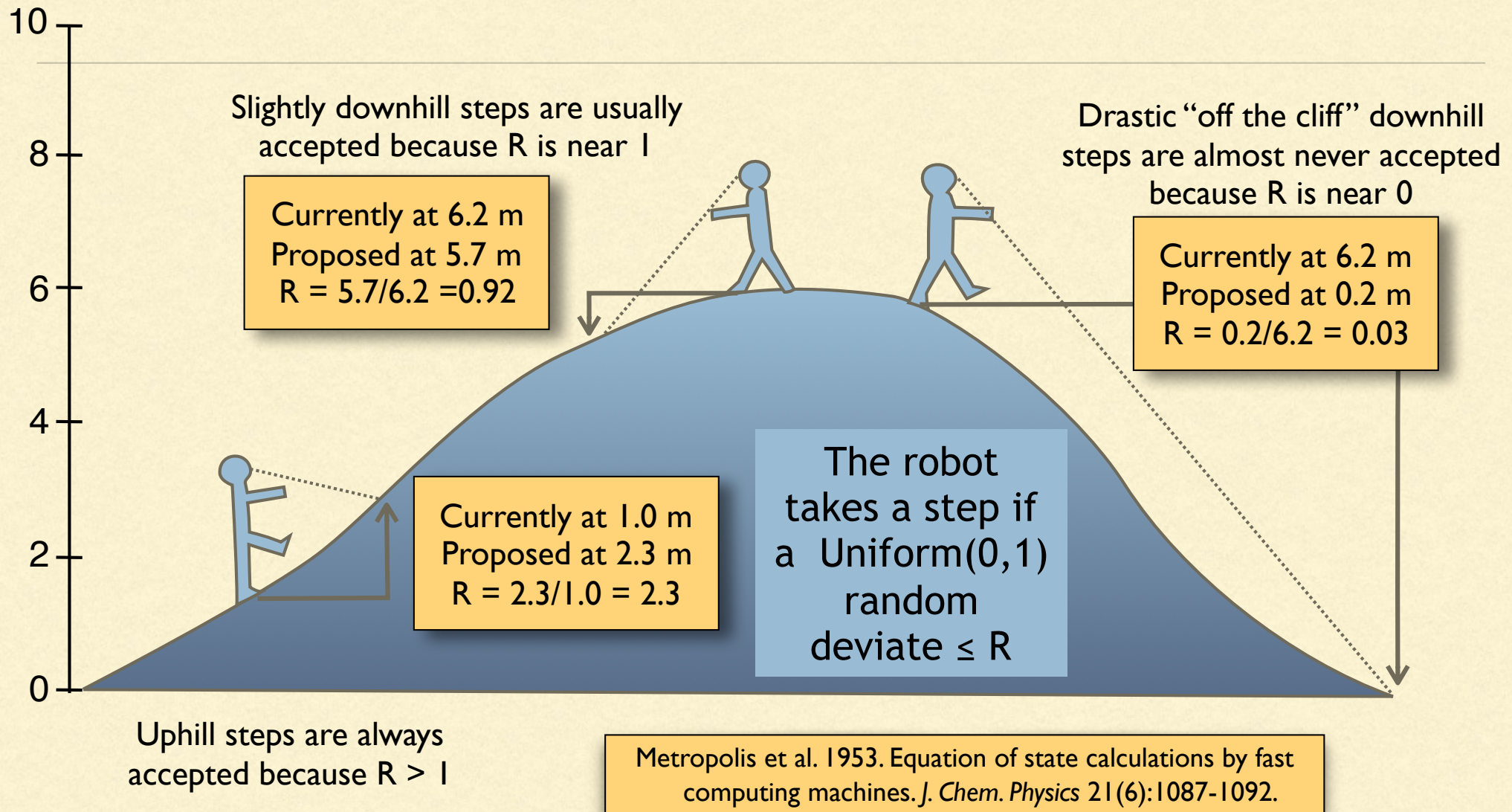


For more complex problems,
we might settle for a
good approximation
to the posterior distribution

MCMC robot's rules



Actual rules (Metropolis algorithm)



Cancellation of marginal likelihood

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

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

Posterior
odds

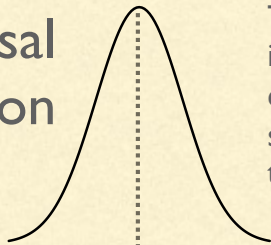
Apply Bayes' rule to
both top and bottom

Likelihood
ratio

Prior
odds

Target vs. Proposal Distributions

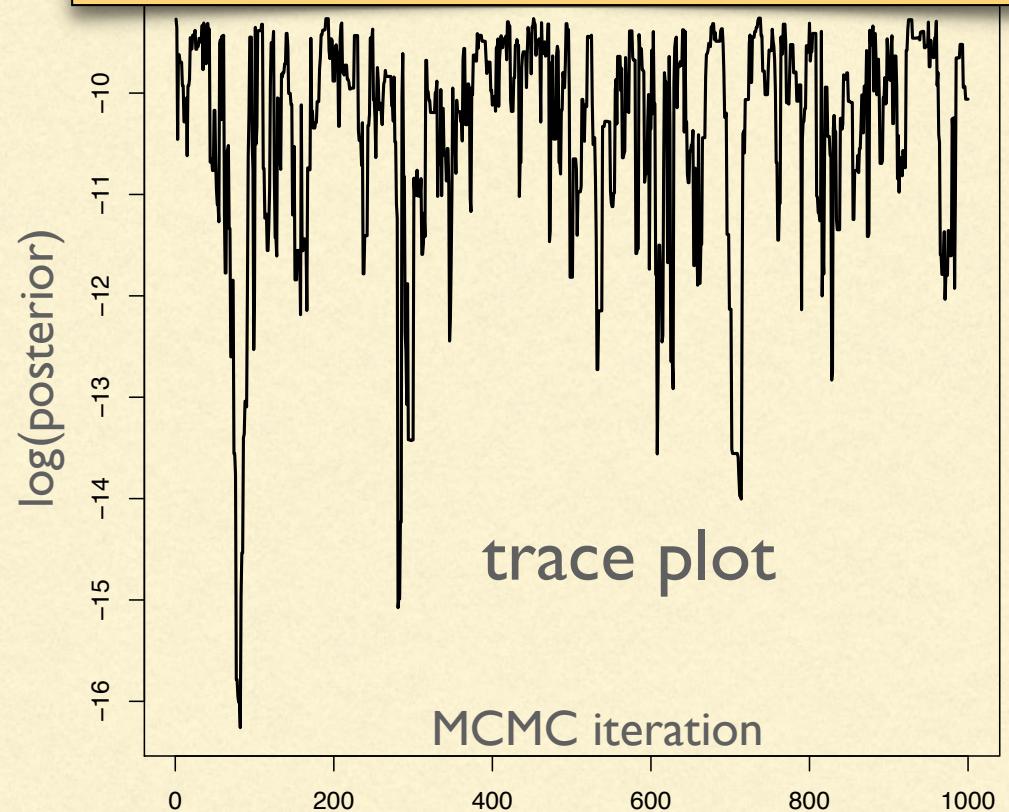
"good" proposal distribution



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

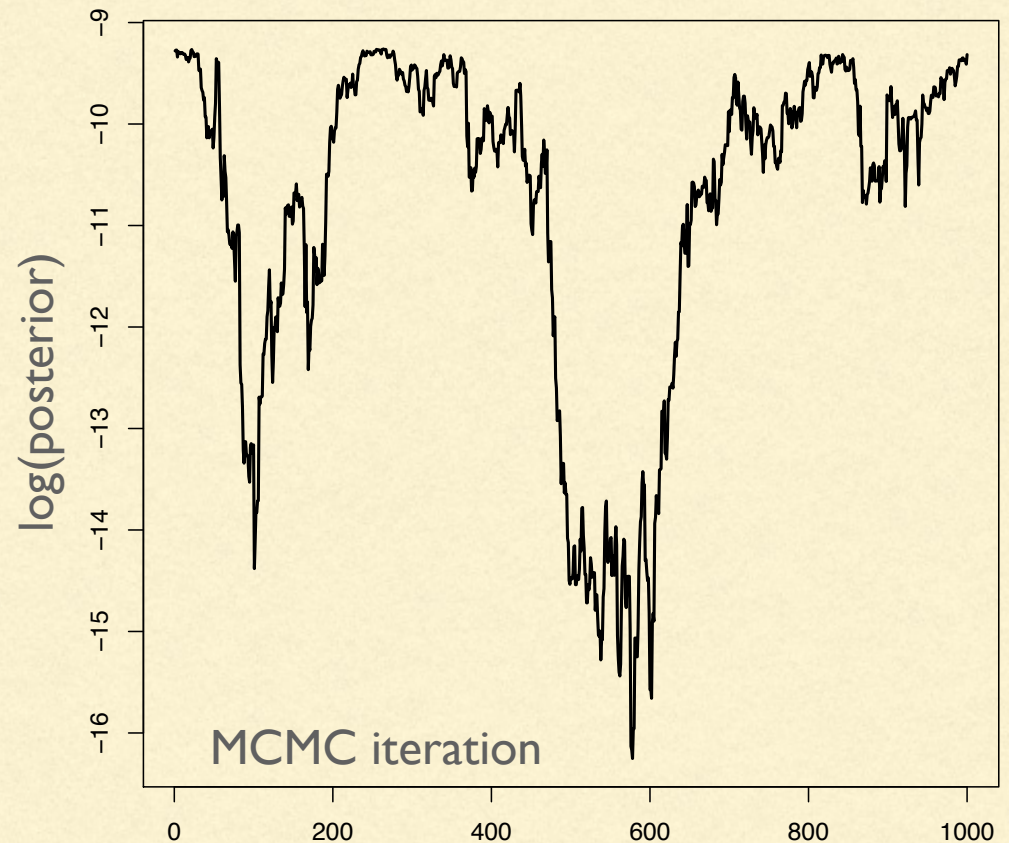
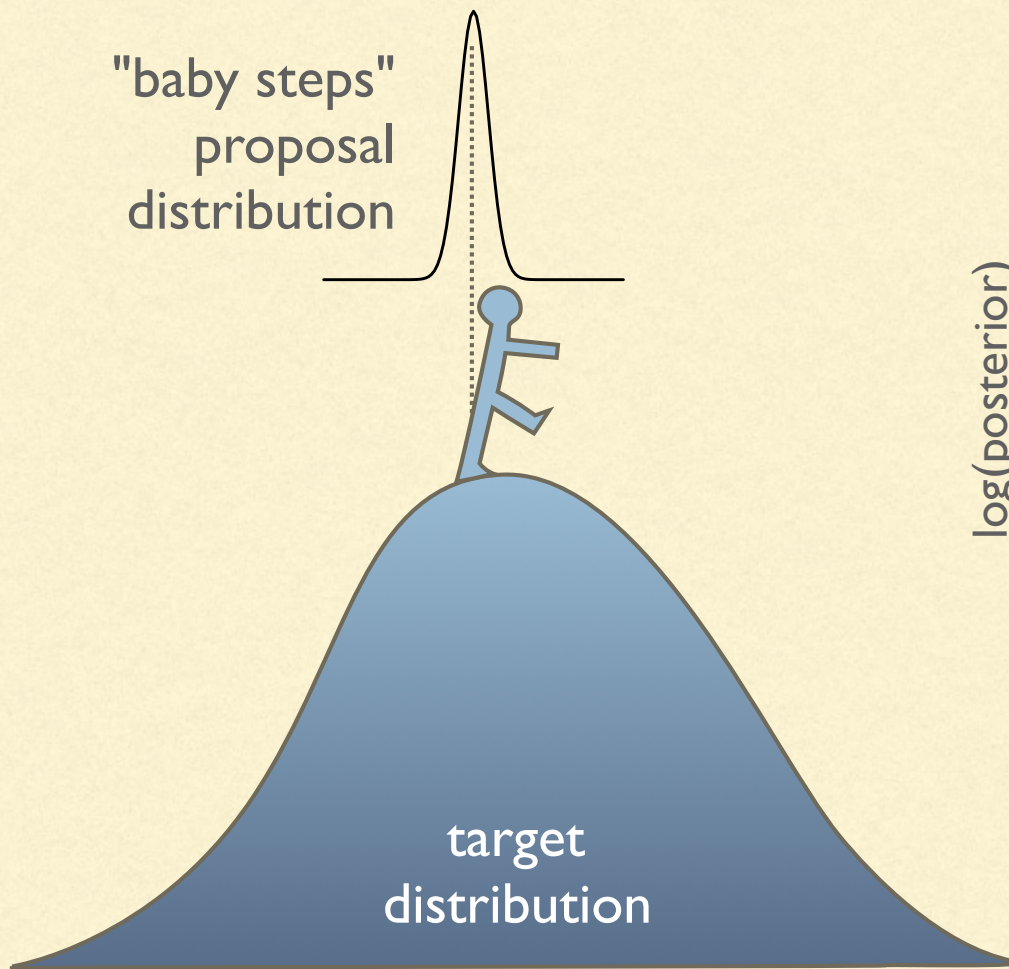


Tracer (app for generating trace plots from MCMC output):
<http://tree.bio.ed.ac.uk/software/tracer/>



White noise appearance is a sign of **good mixing**

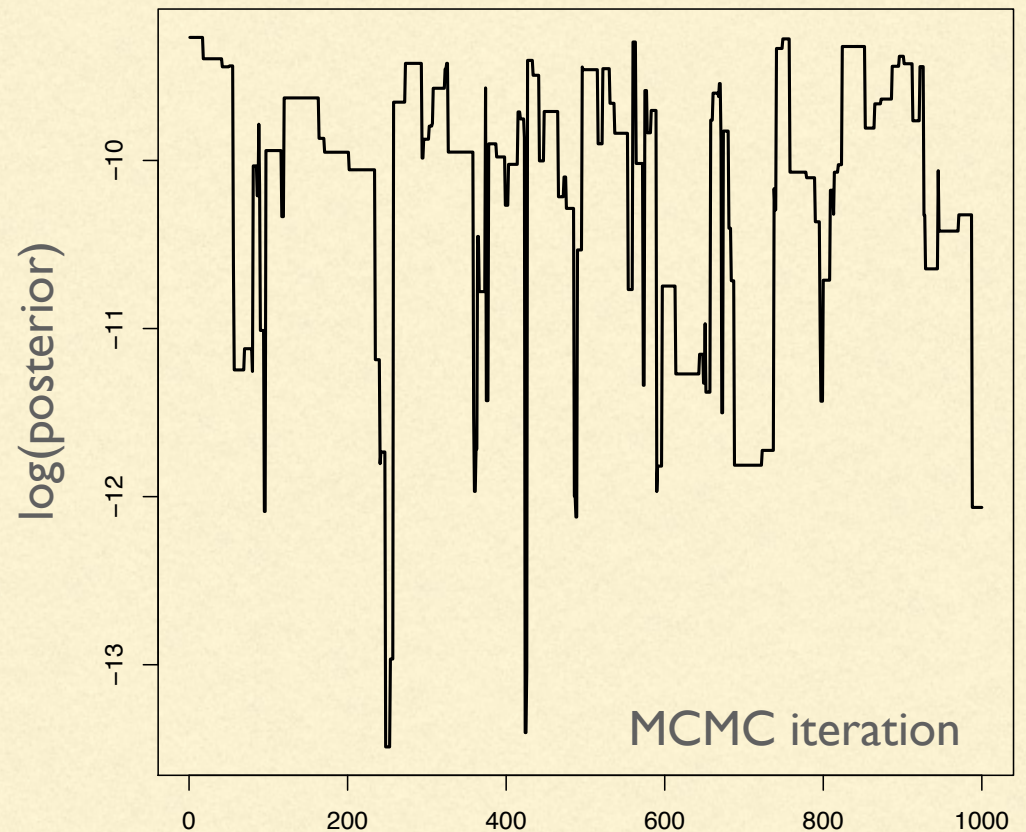
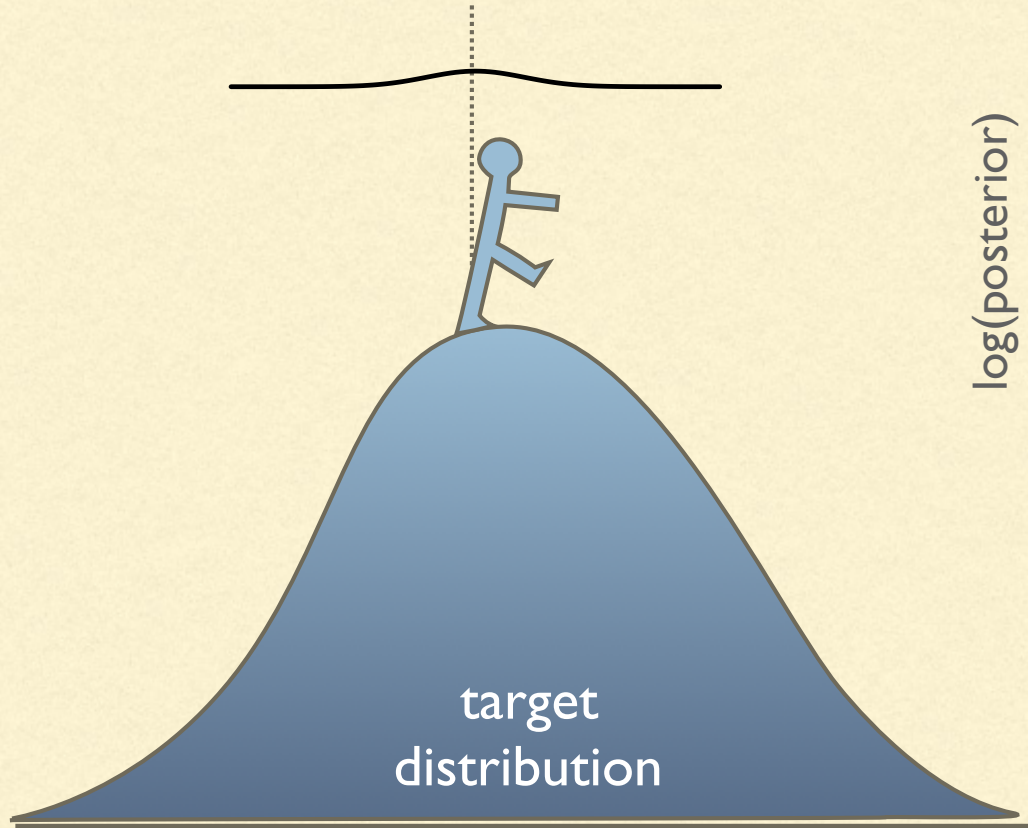
Target vs. Proposal Distributions



Big waves in trace plot indicate robot is crawling around

Target vs. Proposal Distributions

"overly bold" proposal distribution



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:

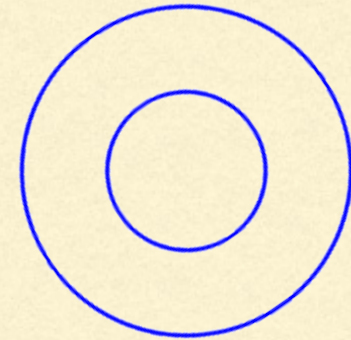
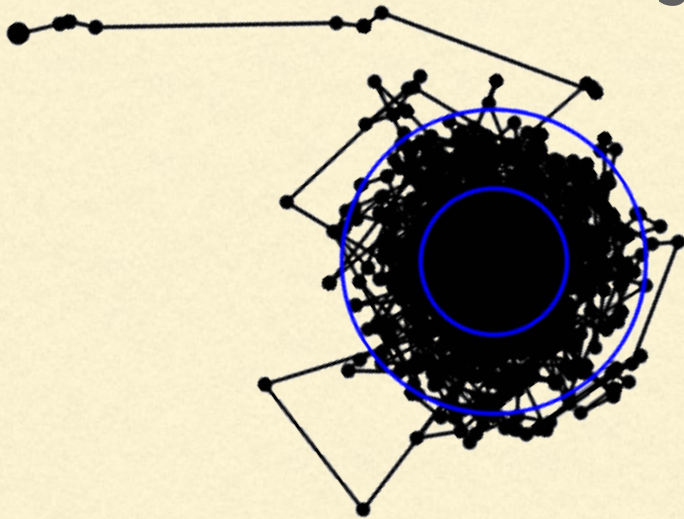
<https://phylogeny.uconn.edu/mcmc-robot/>

Free app for **Windows** or **iPhone/iPad** available
from <http://mcmicrobot.org/>

(also see John Huelsenbeck's iMCMC app for MacOS:
<http://cteg.berkeley.edu/software.html>)

Metropolis-coupled Markov chain Monte Carlo (MCMCMC)

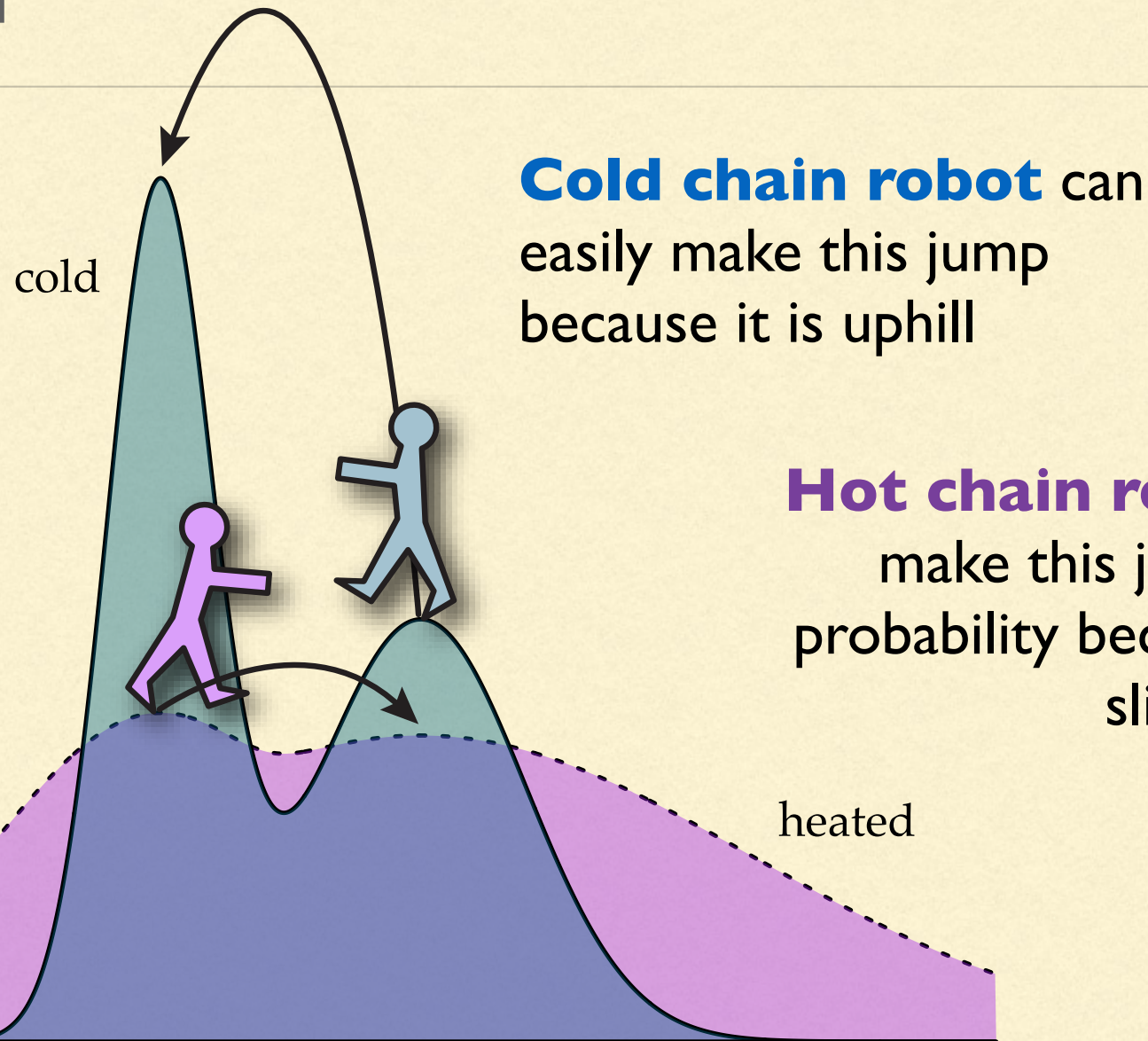
Sometimes the robot needs some help,



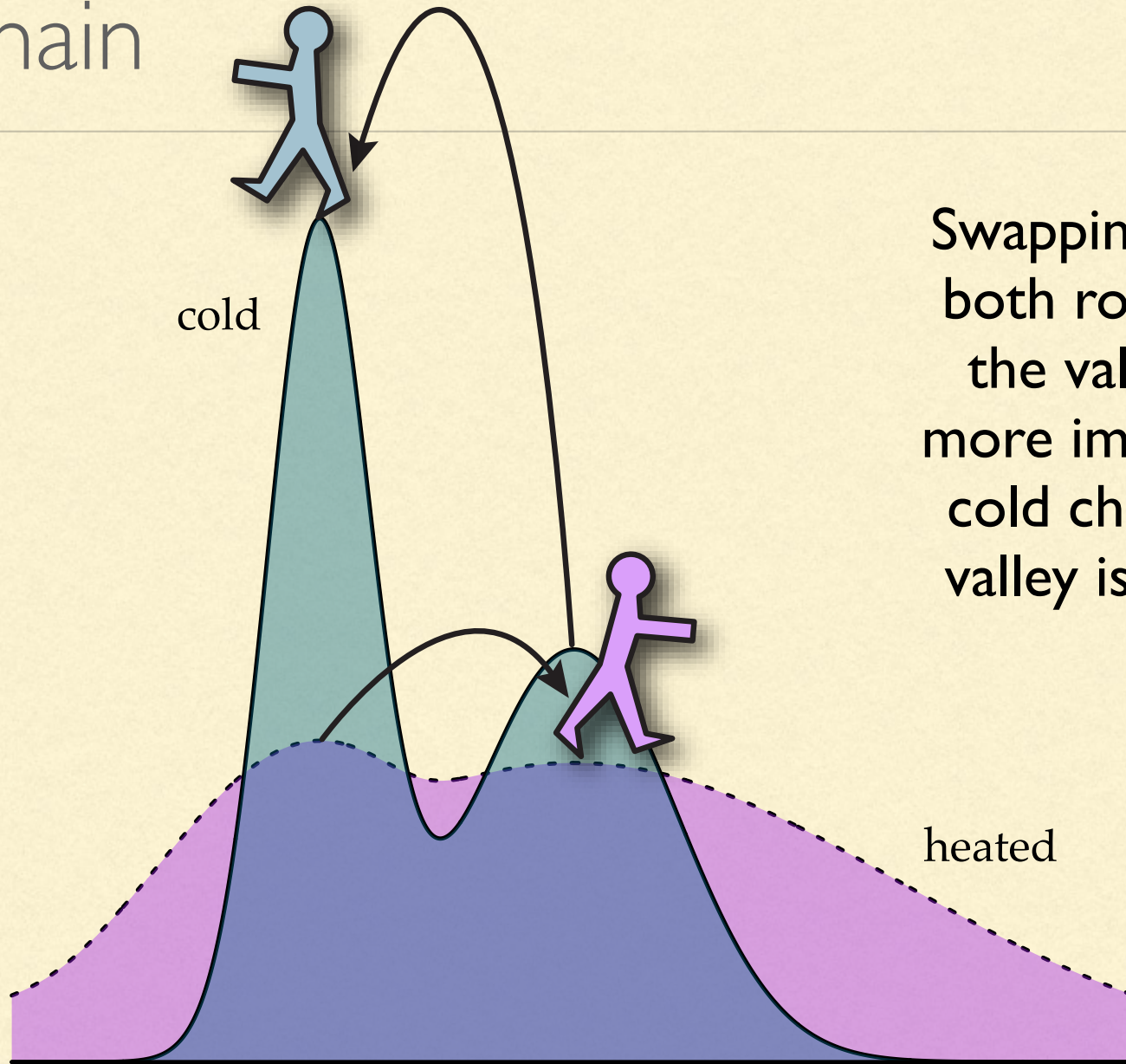
MCMCMC introduces helpers in the form of "heated chain" robots that can act as scouts.

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

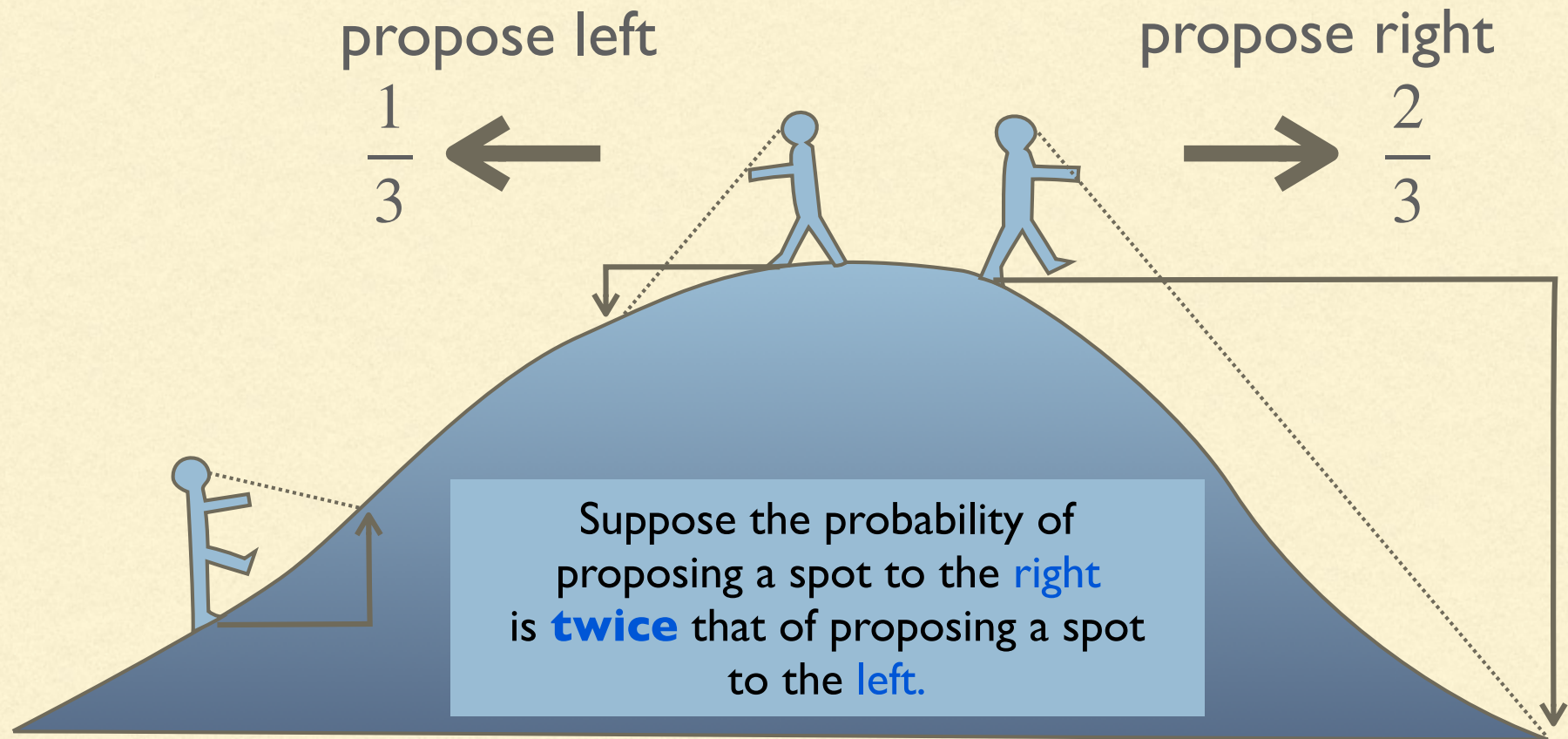


Heated chains act as scouts for the cold chain



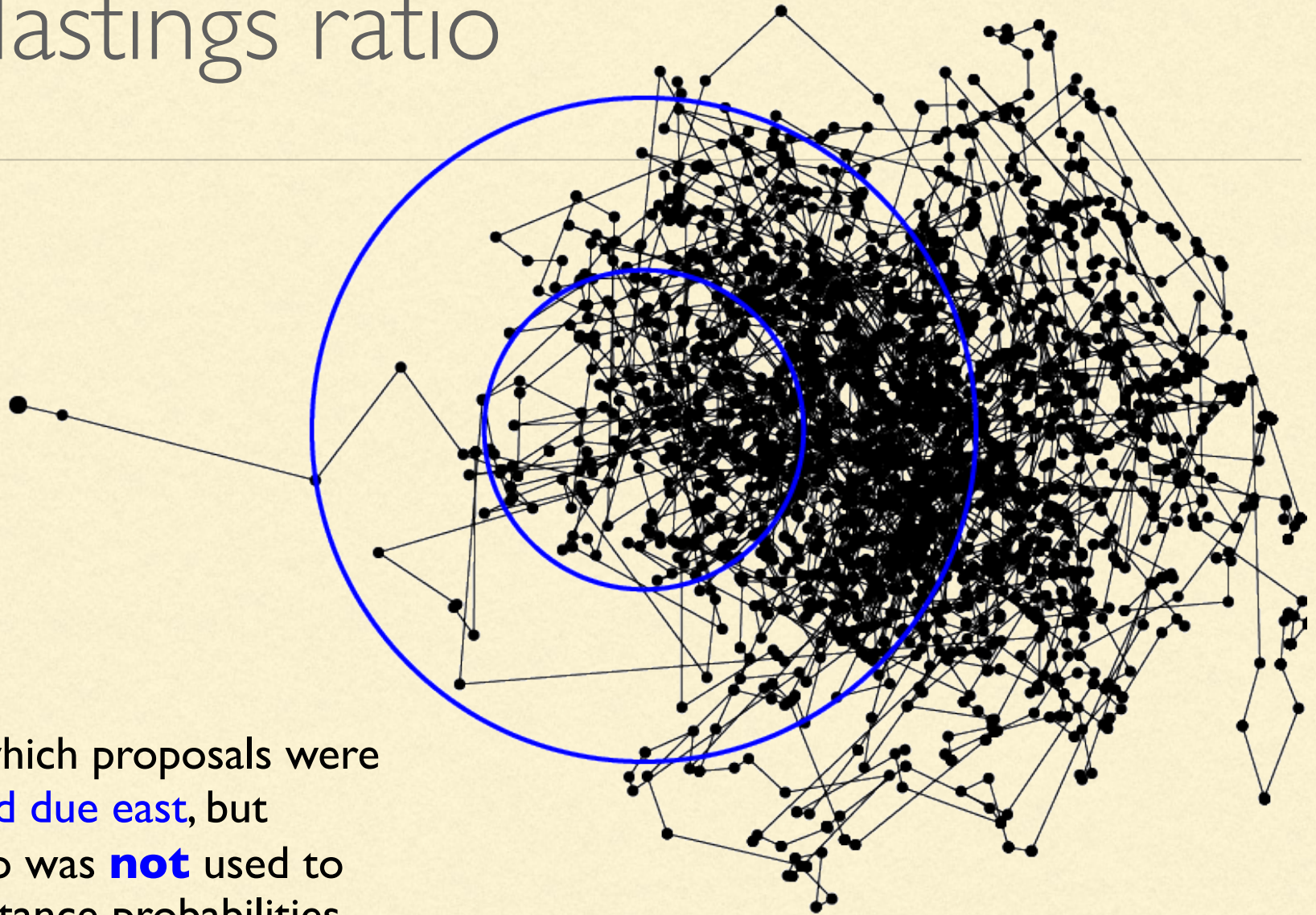
Swapping places means both robots can cross the valley, but this is more important for the cold chain because its valley is much deeper.

The Hastings ratio



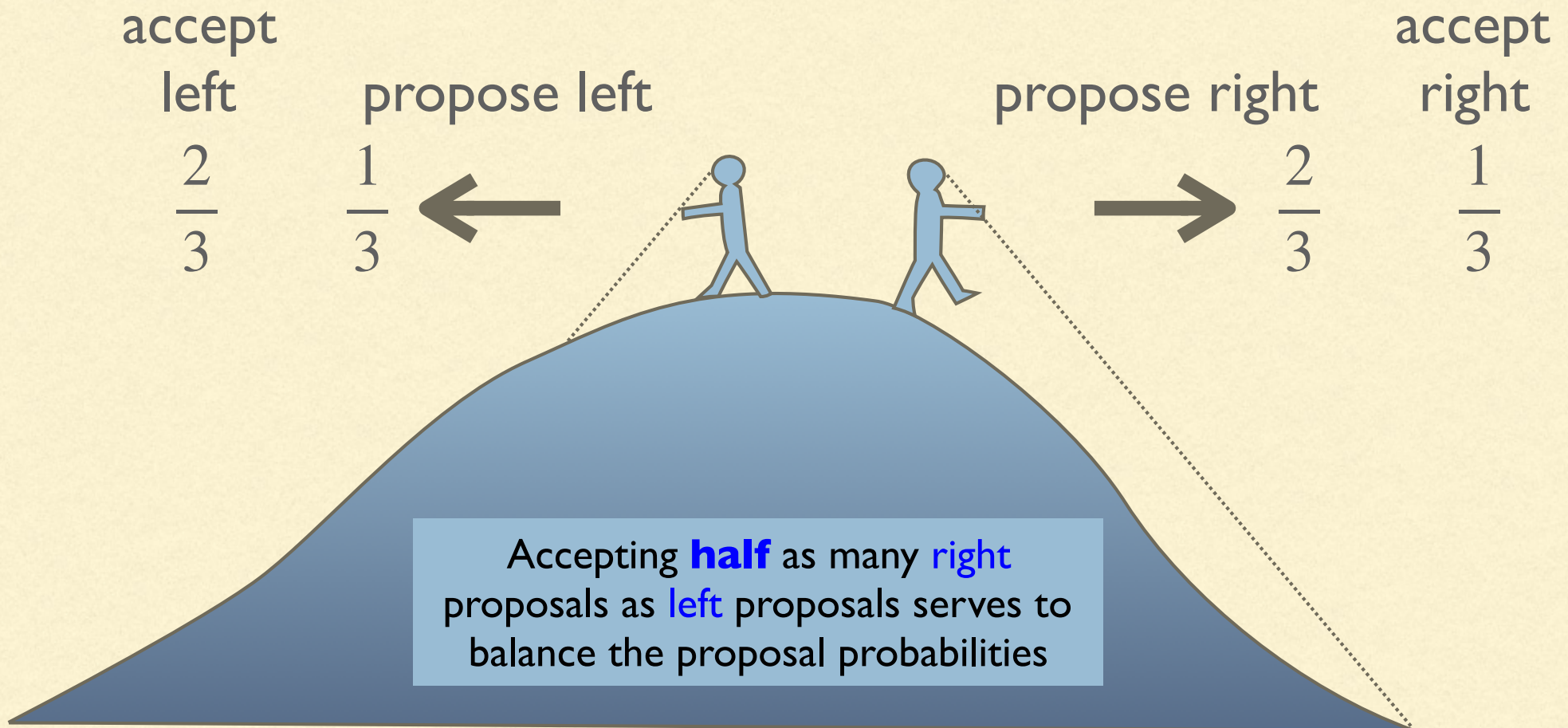
Hastings, W. K. 1970. Monte Carlo sampling methods using Markov chains and their applications. *Biometrika* 57:97-109.

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



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

Hastings Ratio

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

posterior ratio Hastings ratio

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

End of Part 3a

Upcoming Part 3b topics: MCMC proposals, prior distributions, hierarchical models, and Bayes factors