



Introduction to Bayesian Inference

Video Lectures

Paul Lewis's Primer on Phylogenetics

- Trees & Likelihood
- Substitution Models
- Bayesian Statistics & MCMC
- Bayesian Phylogenetics



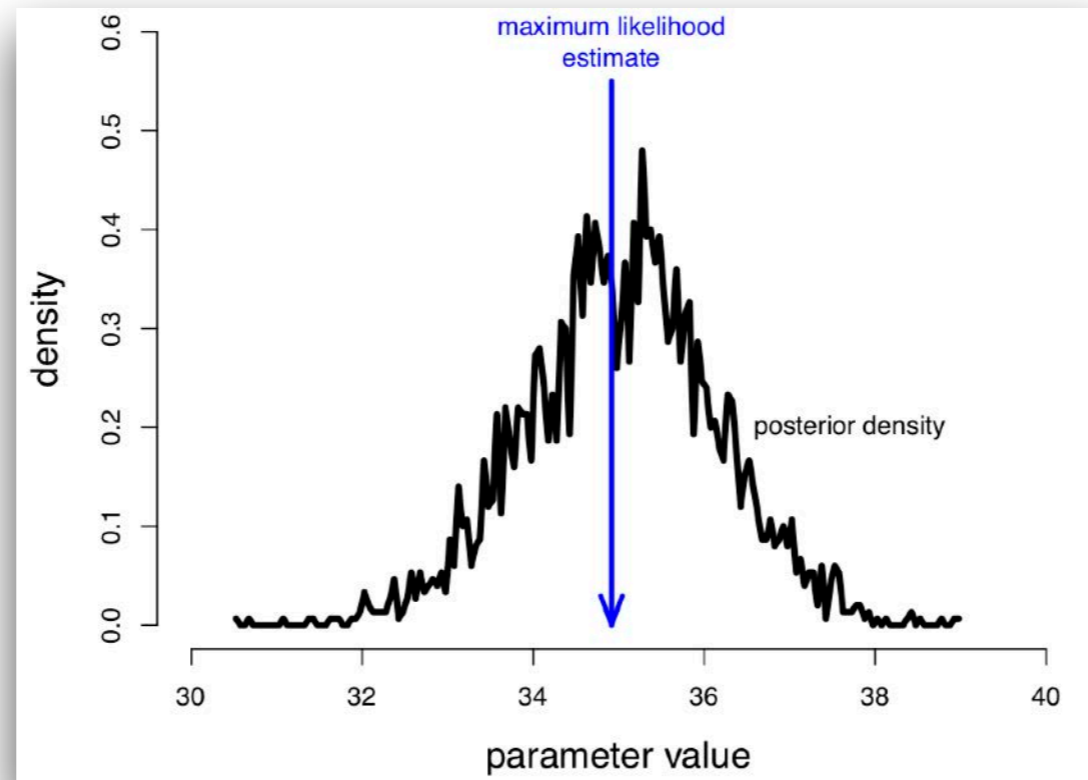
Bayesian or Maximum Likelihood?

Bayesian

- estimates $\Pr(\theta | \mathbf{X})$
- estimates a **distribution**
- parameters are random variables
- average over nuisance parameters

Likelihood Maximum

- estimates $\Pr(\mathbf{X} | \theta)$
- **point** estimate
- parameters are fixed/
unknown
- optimize nuisance parameters



Bayes Rule

posterior probability

likelihood

prior probability

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

marginal probability of the data

The diagram illustrates Bayes' Rule. On the left, the posterior probability $\Pr(\theta | D)$ is shown in a yellow box, with an arrow pointing to it from the label "posterior probability". This is equal to a fraction. The numerator consists of two terms: the likelihood $\Pr(D | \theta)$ in a blue box, with an arrow pointing to it from the label "likelihood", and the prior probability $\Pr(\theta)$ in a pink box, with an arrow pointing to it from the label "prior probability". The denominator is the sum over all θ of the product of the likelihood and prior probability, $\sum_{\theta} \Pr(D | \theta) \Pr(\theta)$, shown in a green box, with an arrow pointing to it from the label "marginal probability of the data".

Bayesian Inference

Estimate the probability of a hypothesis (model) conditional on observed data

The probability represents a **researcher's degree of belief**

Bayes Rule (also called Bayes Theorem) specifies the conditional probability of the hypothesis given the data

Bayes Rule

the posterior probability of a discrete parameter δ conditional on the data D is

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

the likelihood marginalized over all possible values of δ

Bayes Rule

the posterior probability of a discrete parameter θ conditional on the data D is

$$f(\theta | D) = \frac{f(D | \theta)f(\theta)}{\int_{\theta} f(D | \theta)f(\theta)}$$



the likelihood marginalized over all possible values of θ

Priors

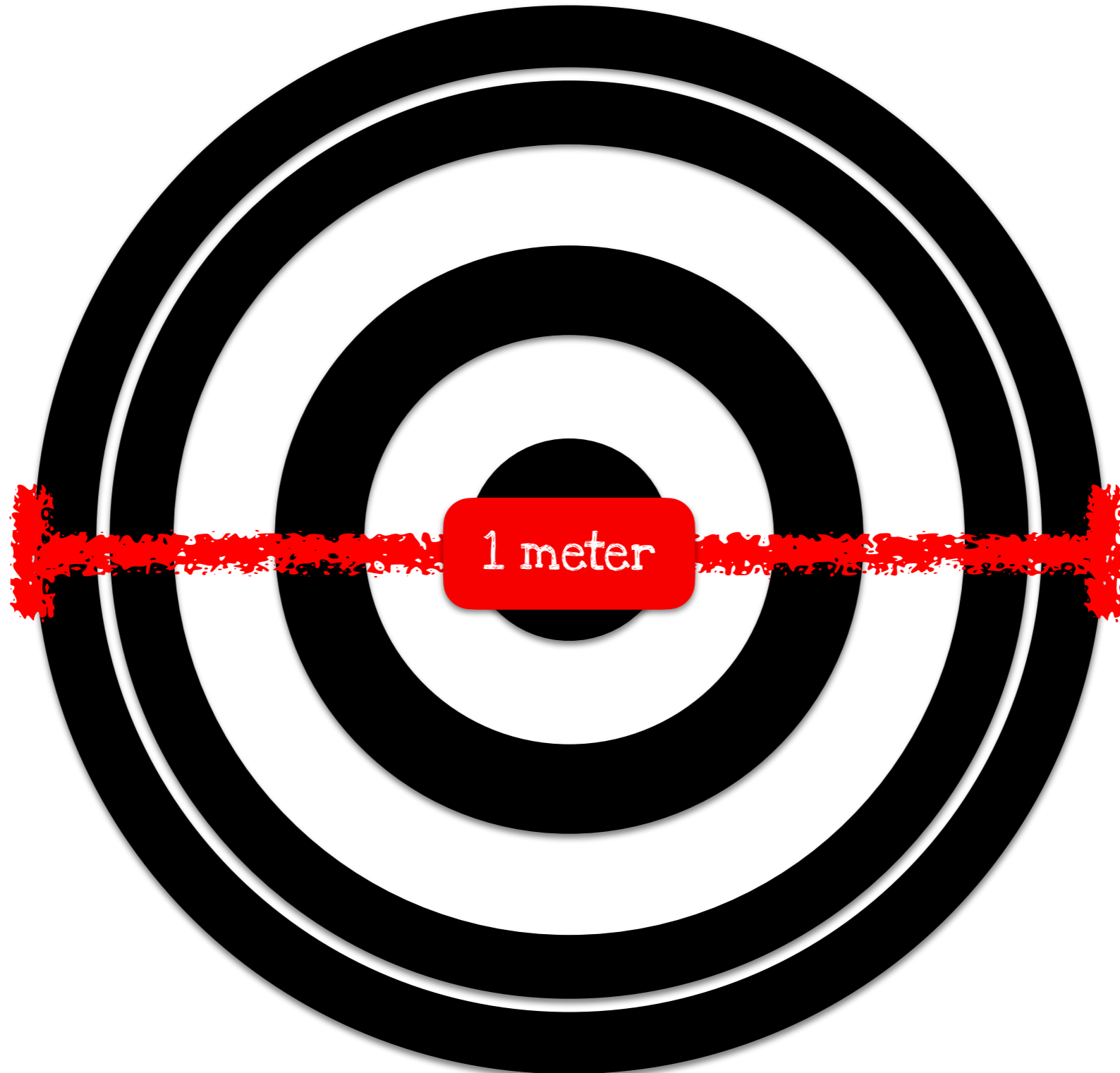
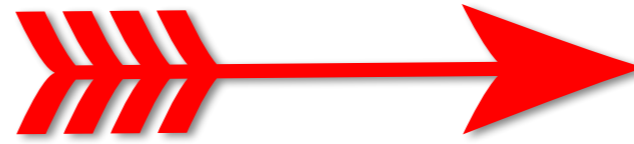
Prior distributions are an important part of Bayesian statistics

The distribution of θ before any data are collected is the prior

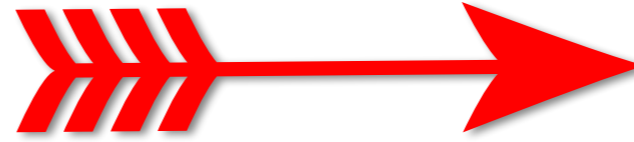
$$f(\theta)$$

The prior describes your uncertainty in the parameters of your model

Priors: Archery



Priors: Archery



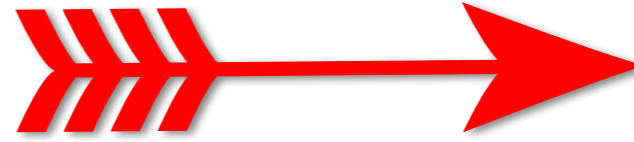
In this example we want to assess an archer's accuracy at hitting the bullseye

To quantify this, we will measure the distance d from the center of the target (in centimeters)



d is an absolute value

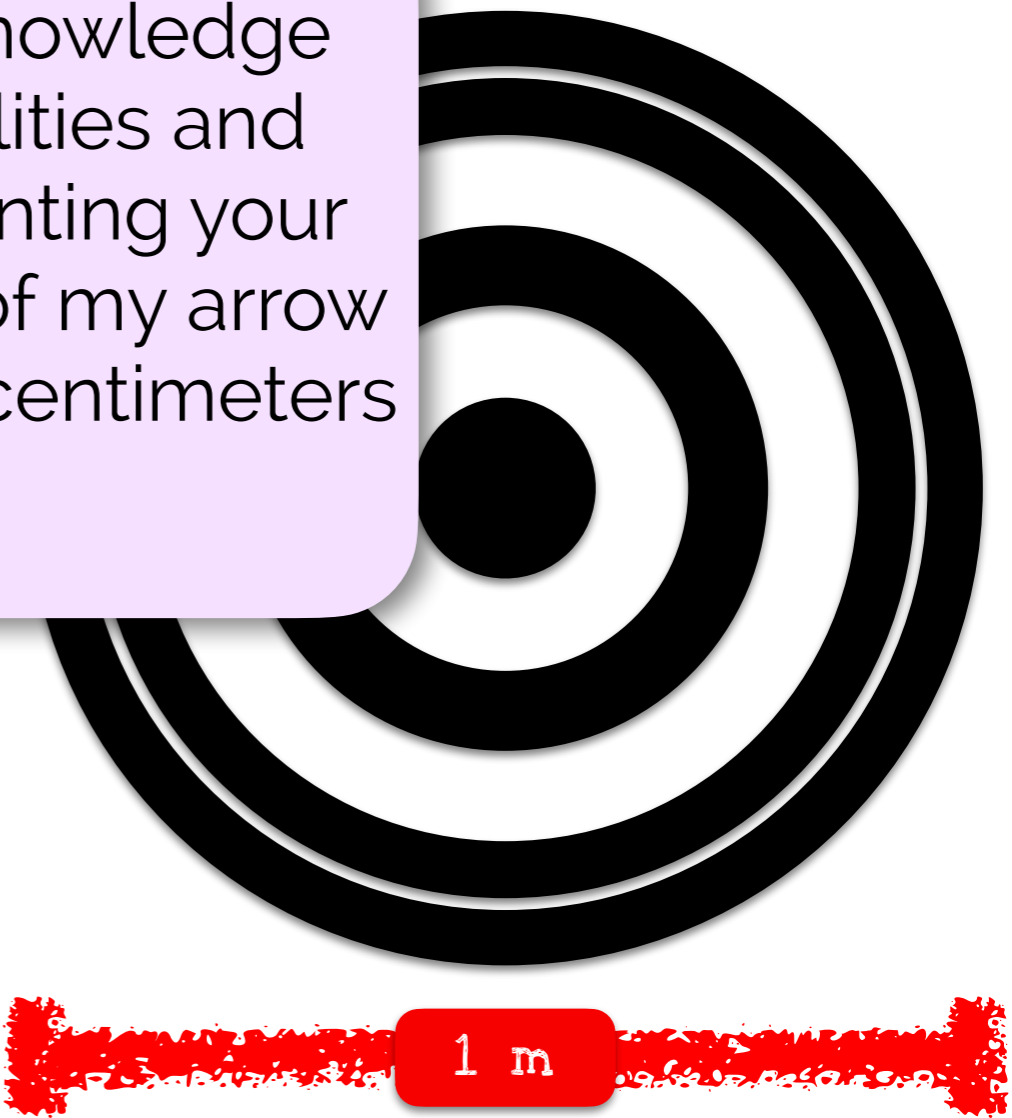
Priors: Archery



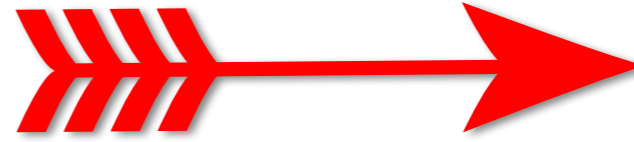
Consider your prior knowledge about my archery abilities and draw a curve representing your view of the chances of my arrow landing a distance d centimeters from the bullseye

When formalizing your prior belief, also consider what you know about d

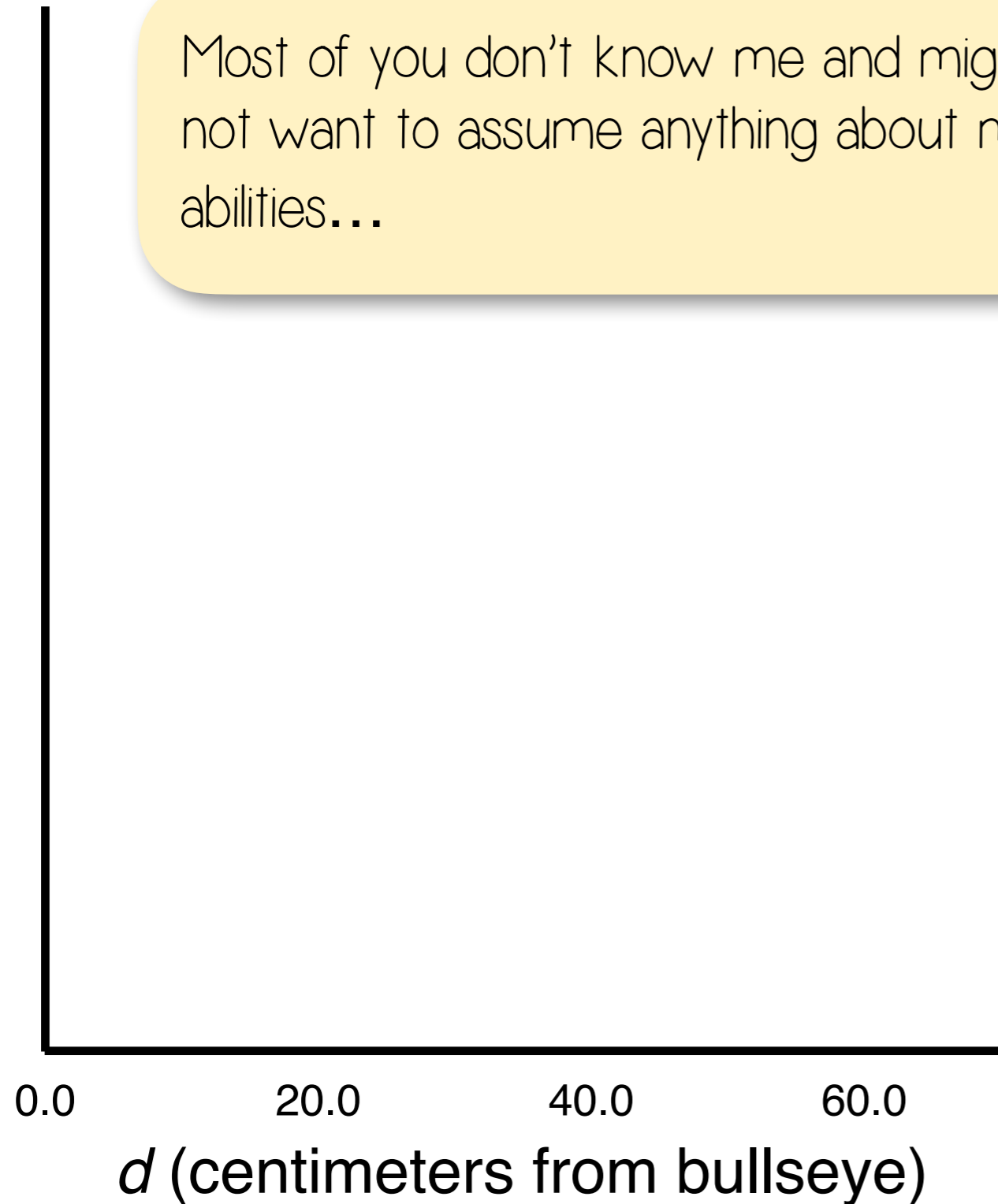
0.0 20.0 40.0 60.0
 d (centimeters from bullseye)



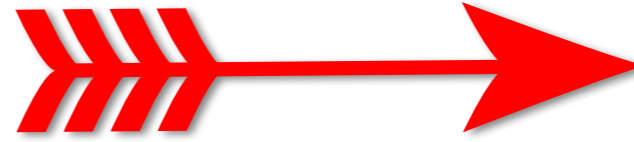
Priors: Archery



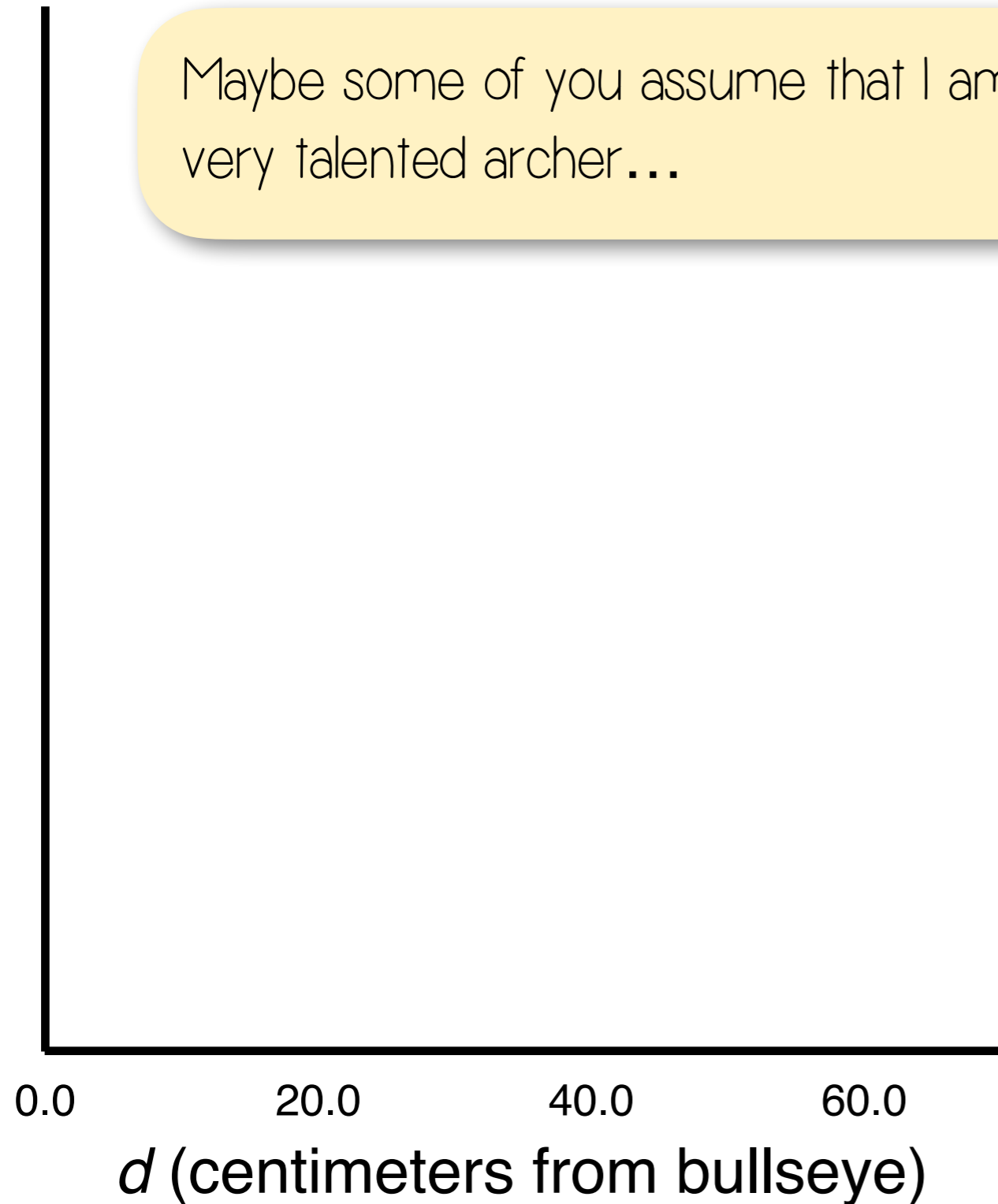
Most of you don't know me and might not want to assume anything about my abilities...



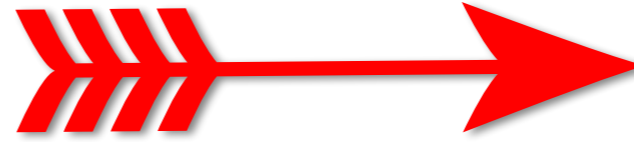
Priors: Archery



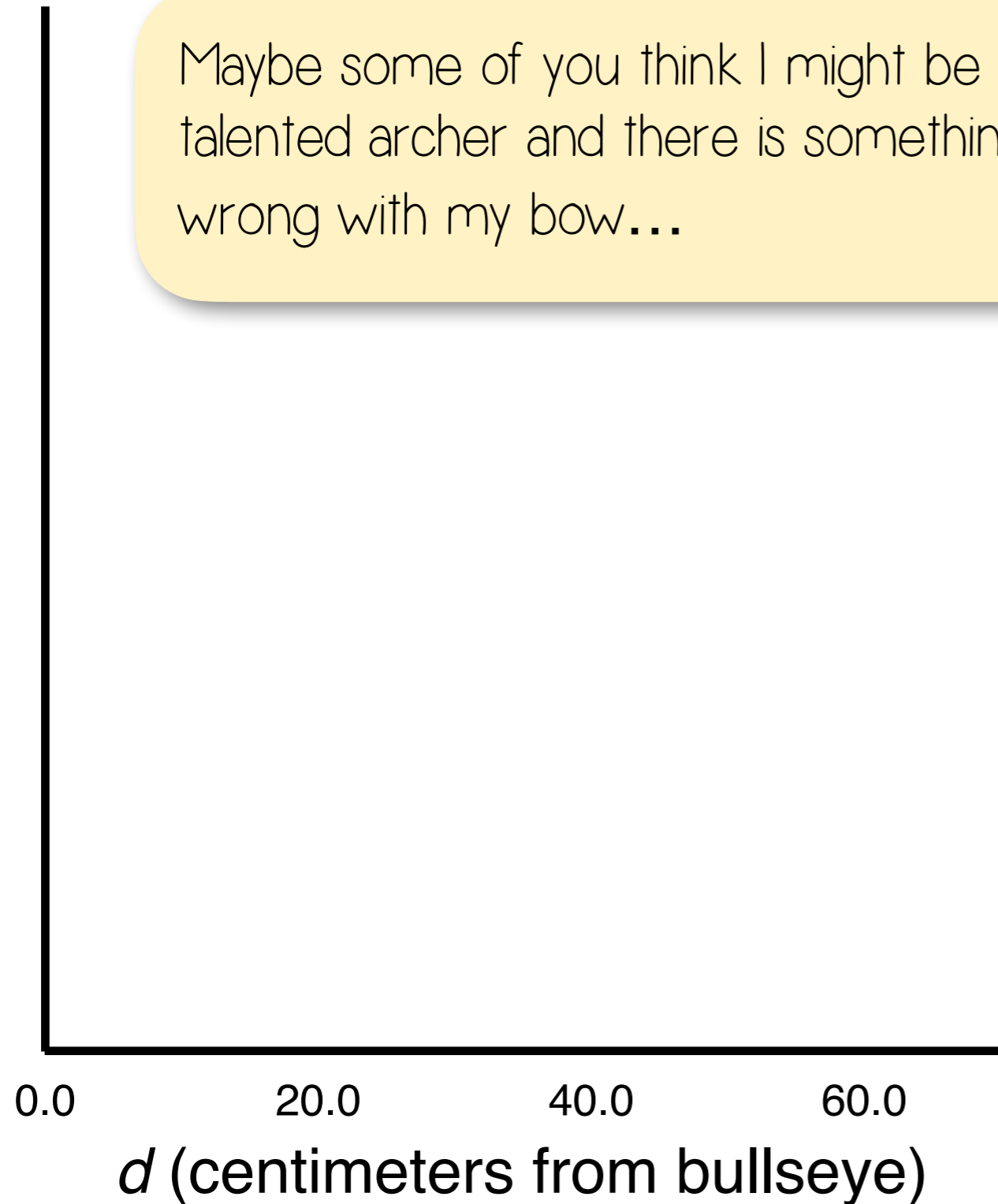
Maybe some of you assume that I am a very talented archer...



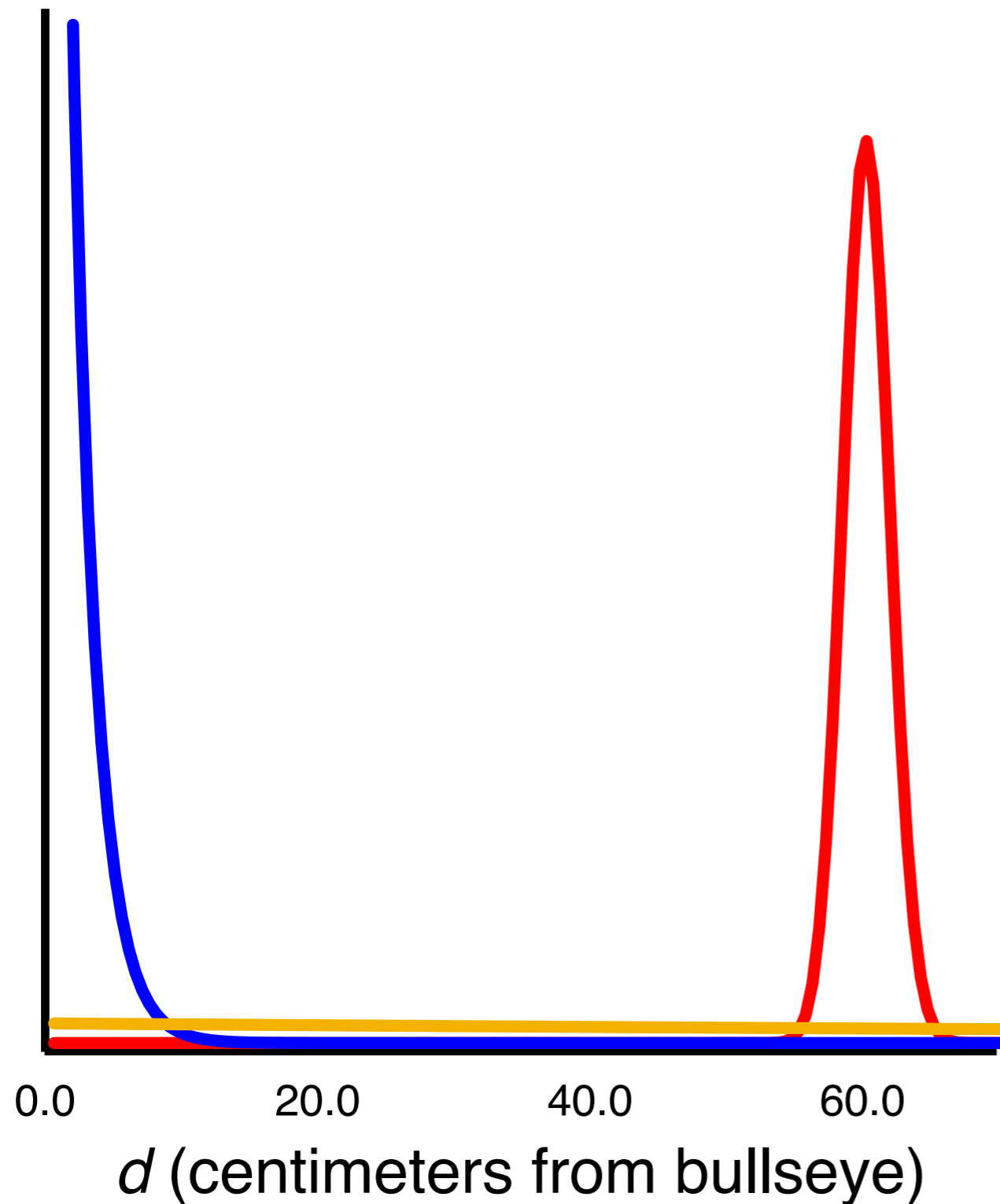
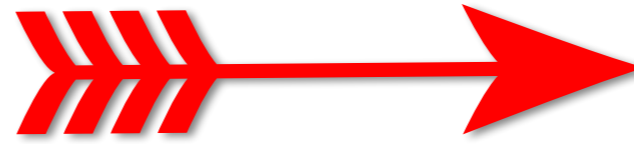
Priors: Archery



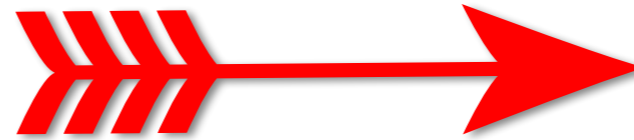
Maybe some of you think I might be a talented archer and there is something wrong with my bow...



Priors: Archery

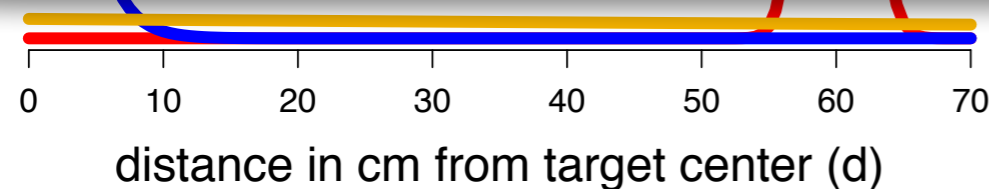


Priors: Archery



Each of these prior densities can be defined using a gamma distribution.

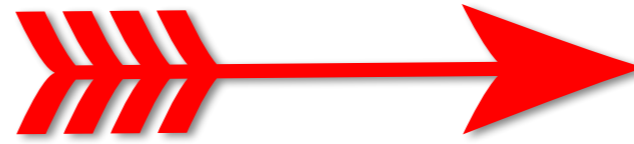
To specify a gamma prior, we must choose parameter values based on our **prior belief**



$$d \sim \text{Gamma}(\alpha, \beta)$$

$$f(d \mid \alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta^\alpha} d^{\alpha-1} e^{-\frac{d}{\beta}}$$

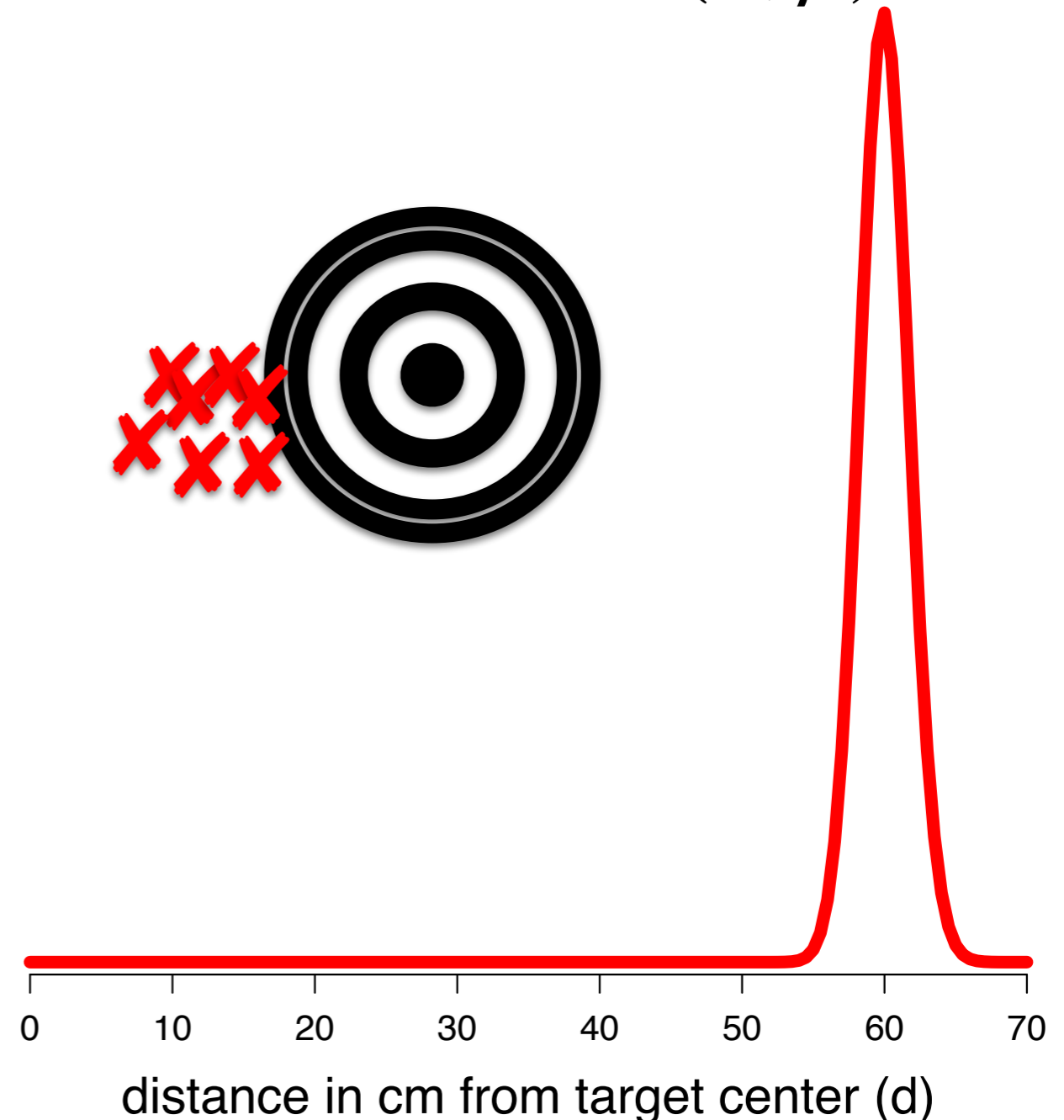
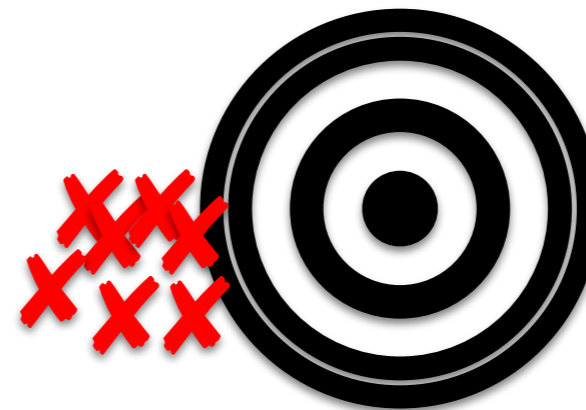
Priors: Archery



Let's assume that I will consistently miss the target

$$d \sim \text{Gamma}(\alpha, \beta)$$

This is a gamma distribution with a mean (m) of 60 and a variance (v) of 3



mean = accuracy

variance = precision

Priors: Archery

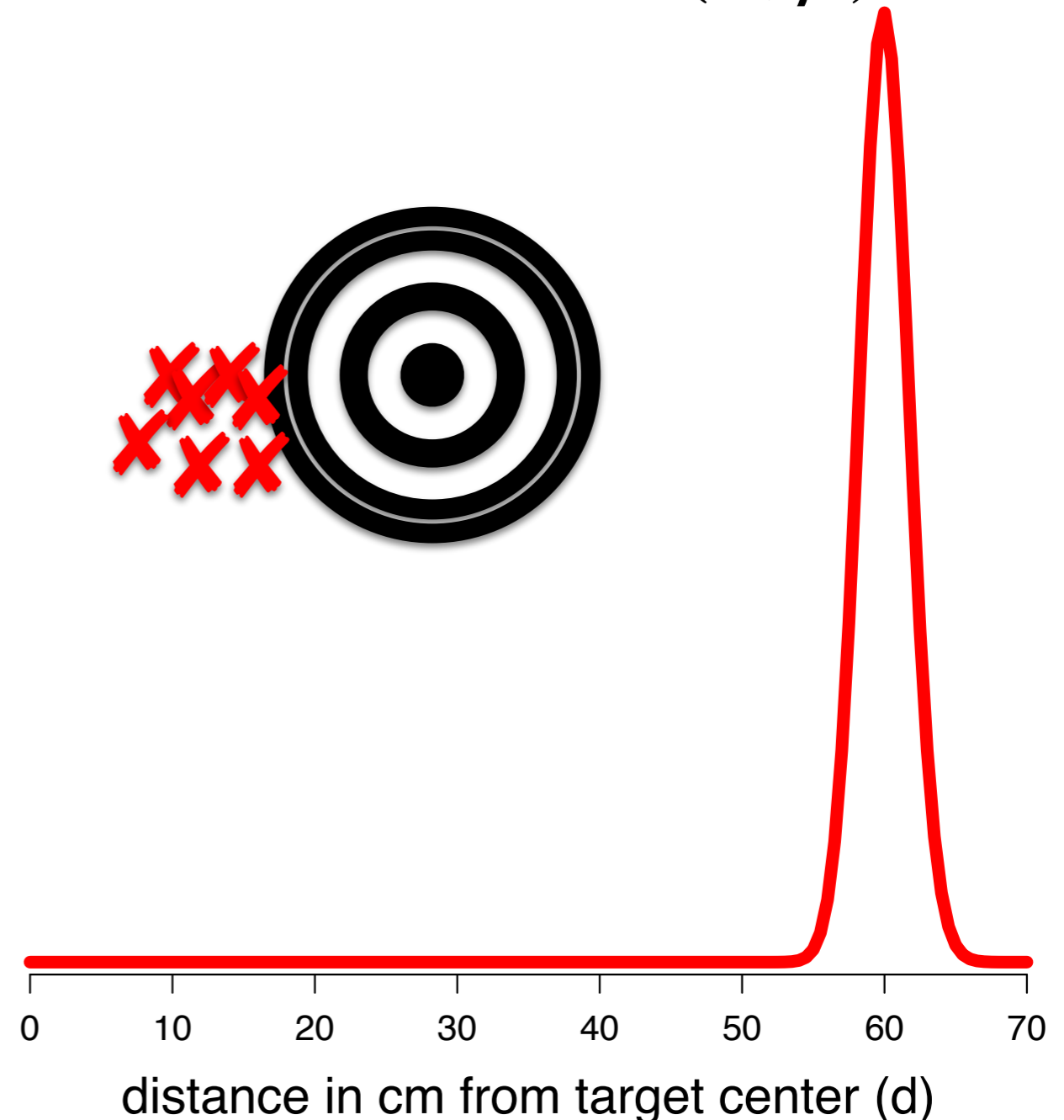
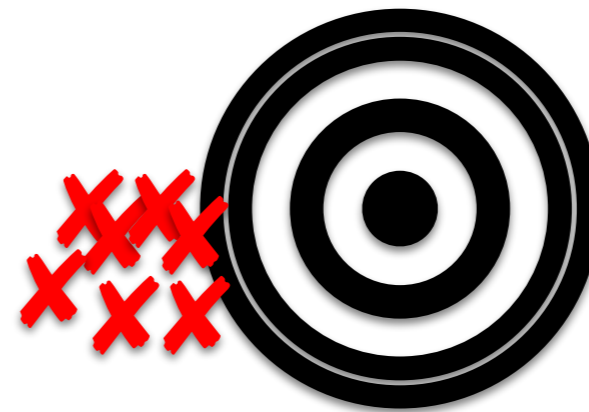


If we have prior knowledge of the mean and variance of the gamma distribution, we can compute the shape and rate parameters

$$m = \frac{\alpha}{\beta}, \quad \alpha = \frac{m^2}{v}$$

$$v = \frac{\alpha}{\beta^2}, \quad \beta = \frac{m}{v}$$

$$d \sim \text{Gamma}(\alpha, \beta)$$



Priors: Archery

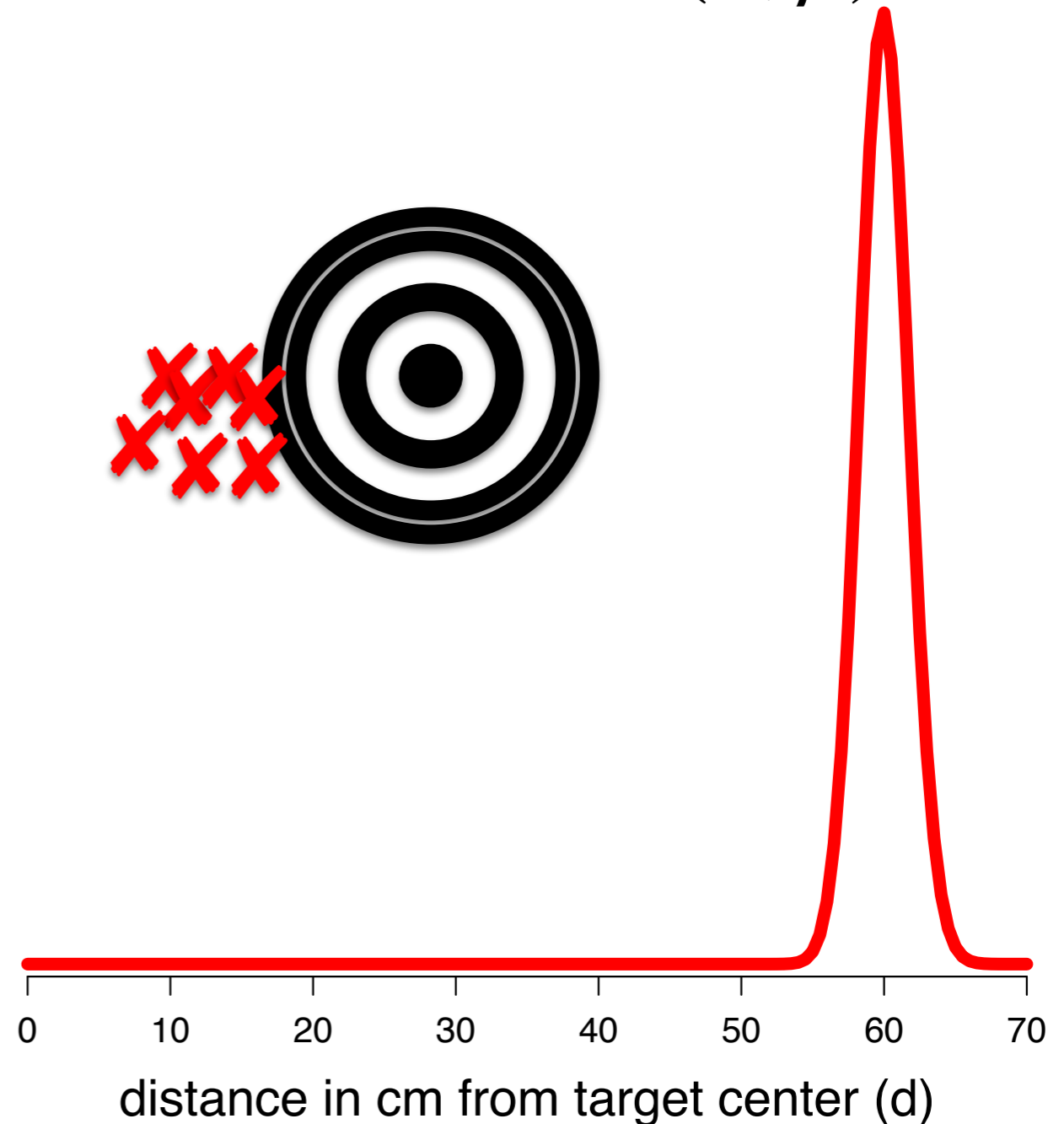
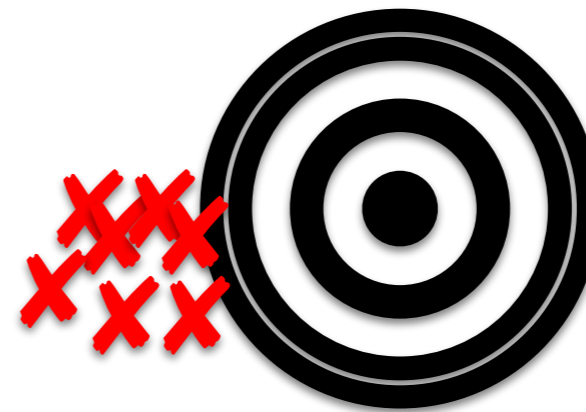


$$m = 60, \quad v = 3$$

$$d \sim \text{Gamma}(\alpha, \beta)$$

$$\alpha = \frac{60^2}{3} = 1200$$

$$\beta = \frac{60}{3} = 20$$

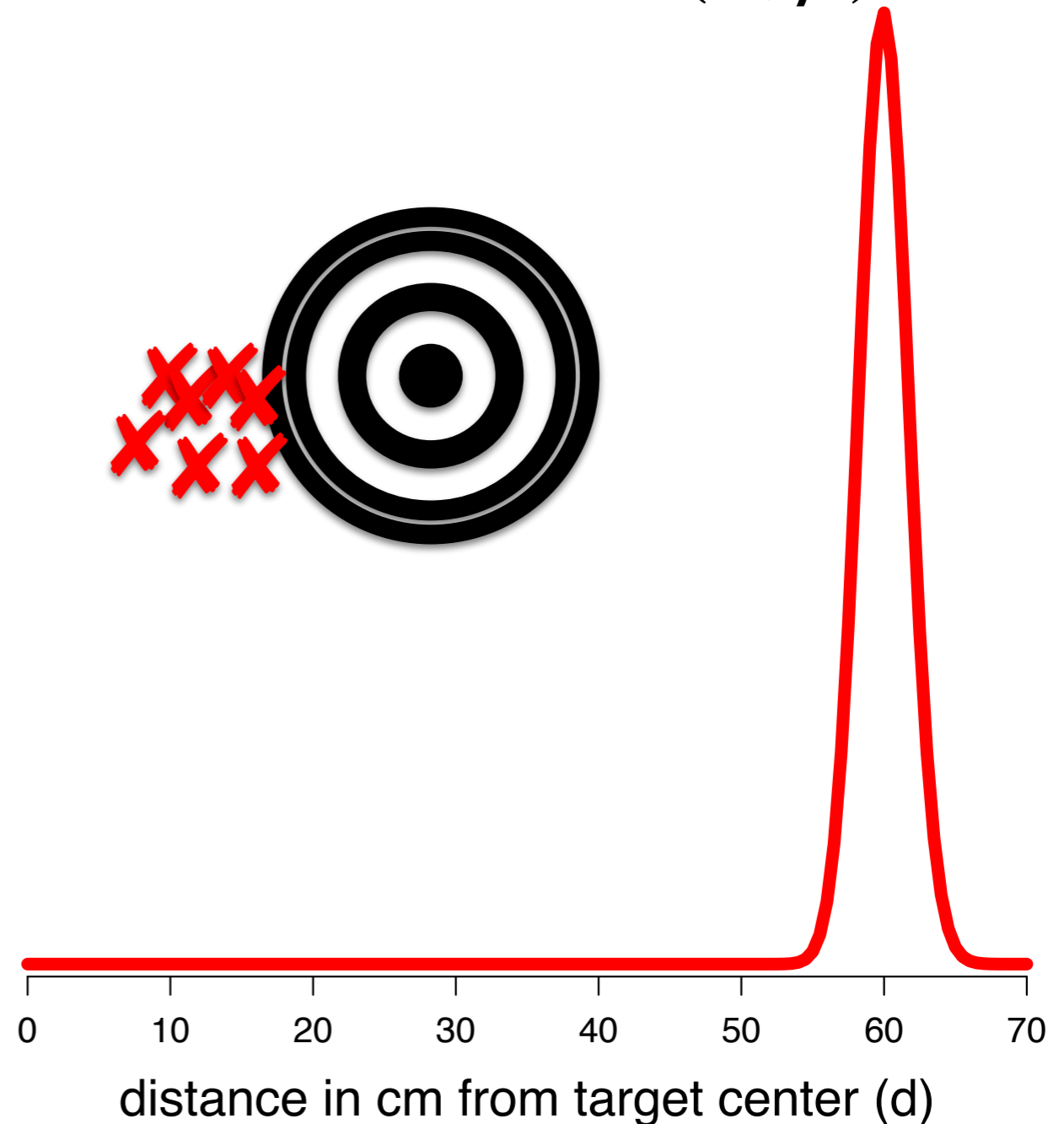
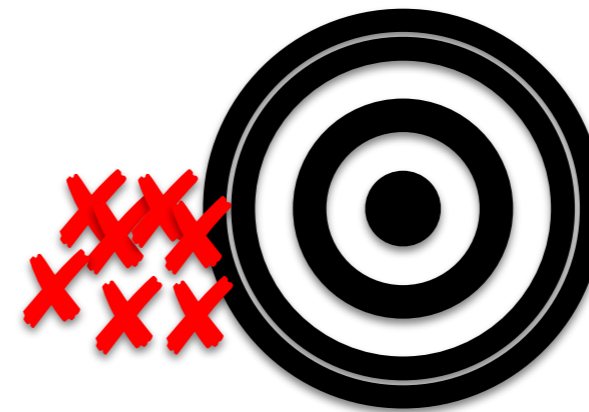
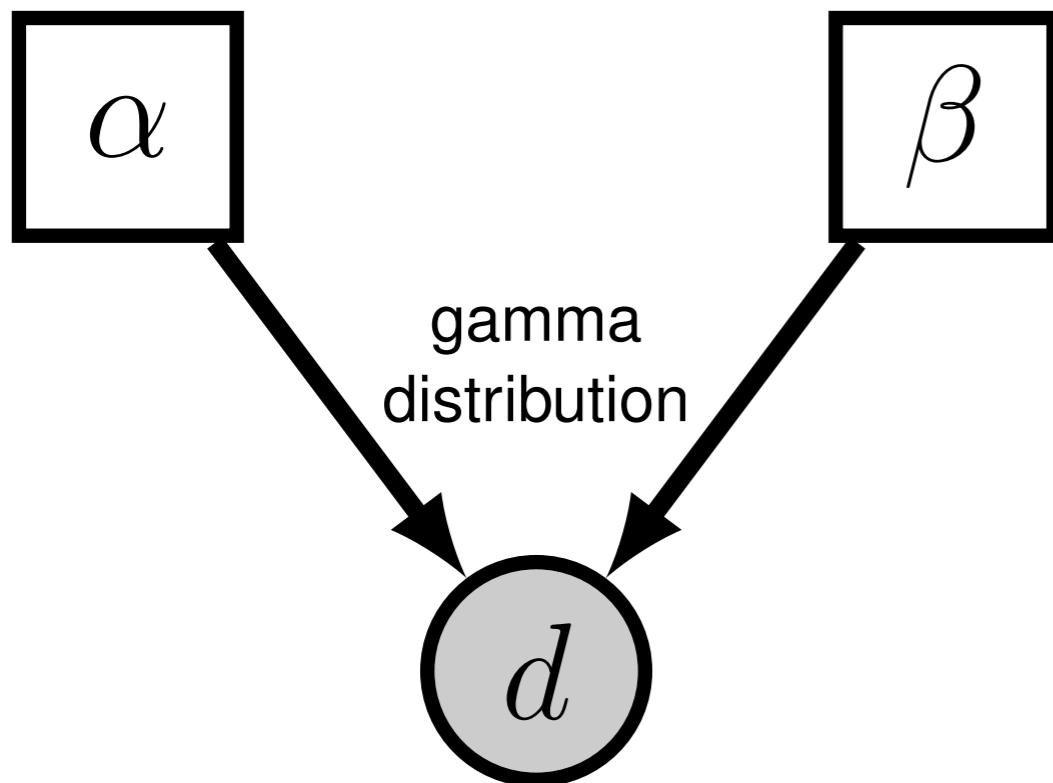


Priors: Archery

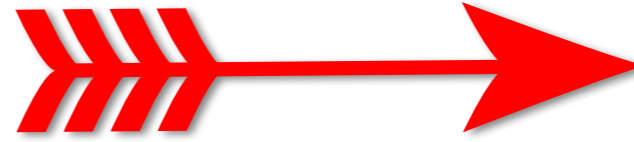


Another way of expressing this distribution is with a probabilistic graphical model

$$d \sim \text{Gamma}(\alpha, \beta)$$

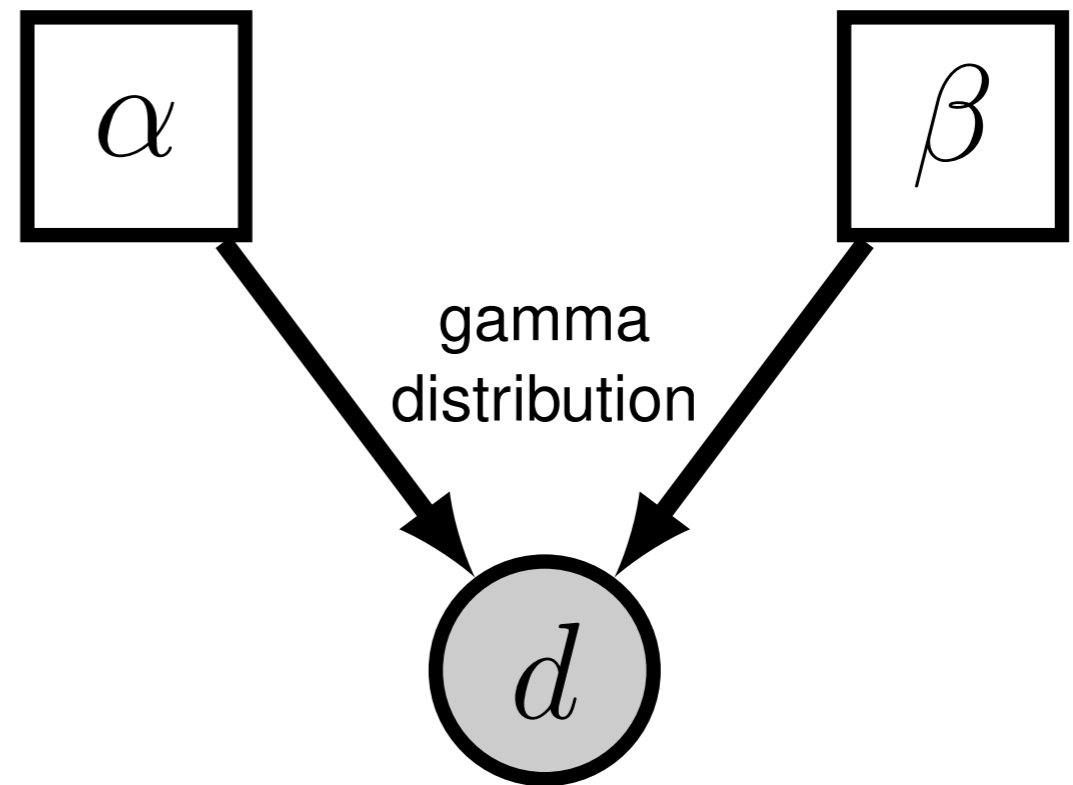


Priors: Archery

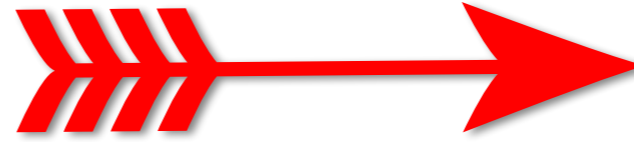


This shows that our observed datum (d = a single observed shot) is conditionally dependent on the shape (α) and rate (β) of the gamma distribution

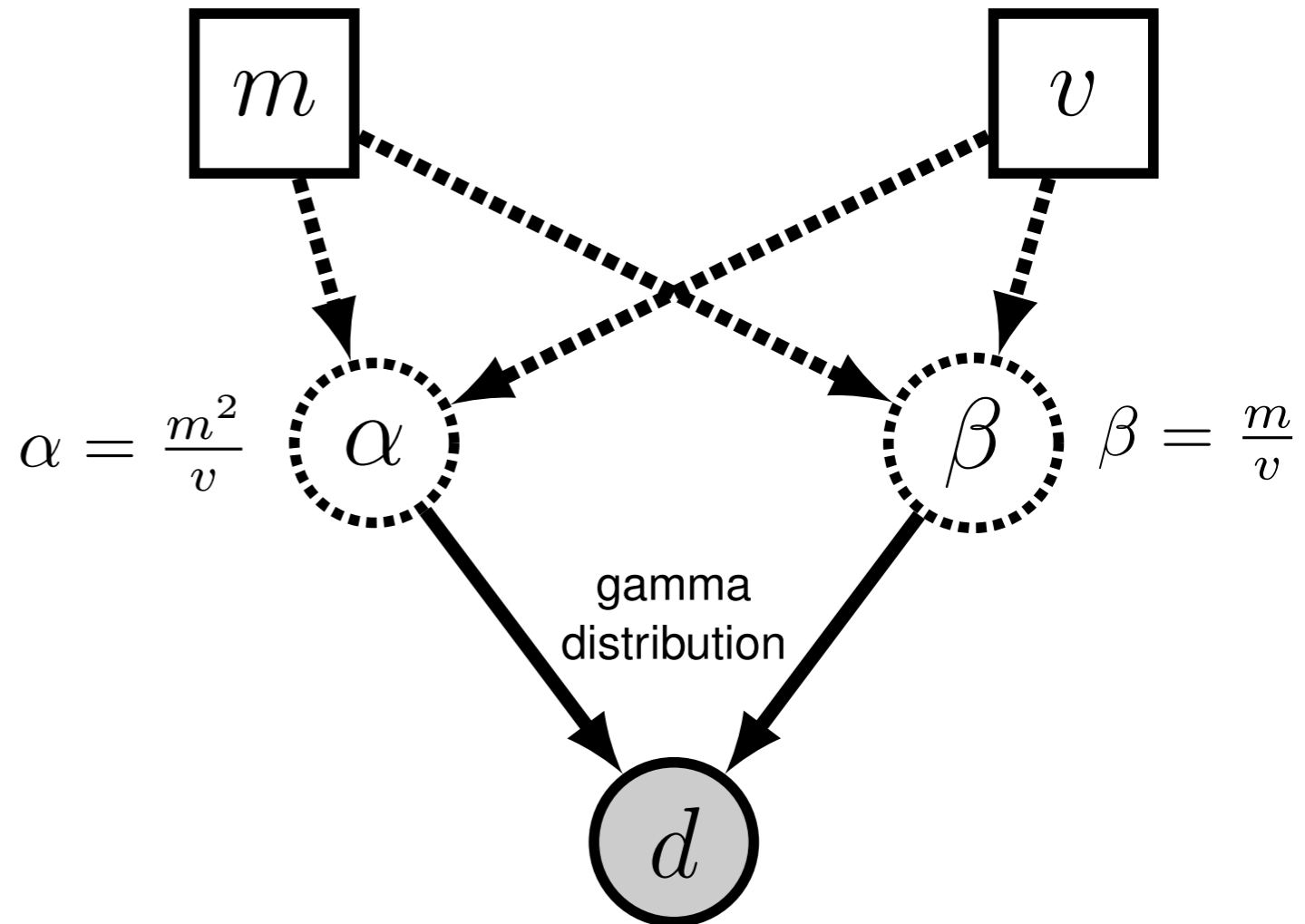
$$d \sim \text{Gamma}(\alpha, \beta)$$



Priors: Archery

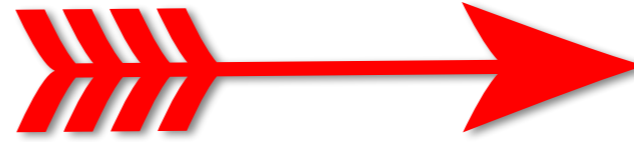


We can parameterize the model using the mean (m) and variance (v), where α and β are computed using m and v



We may have more intuition about the mean and variance than we do about the shape and rate.

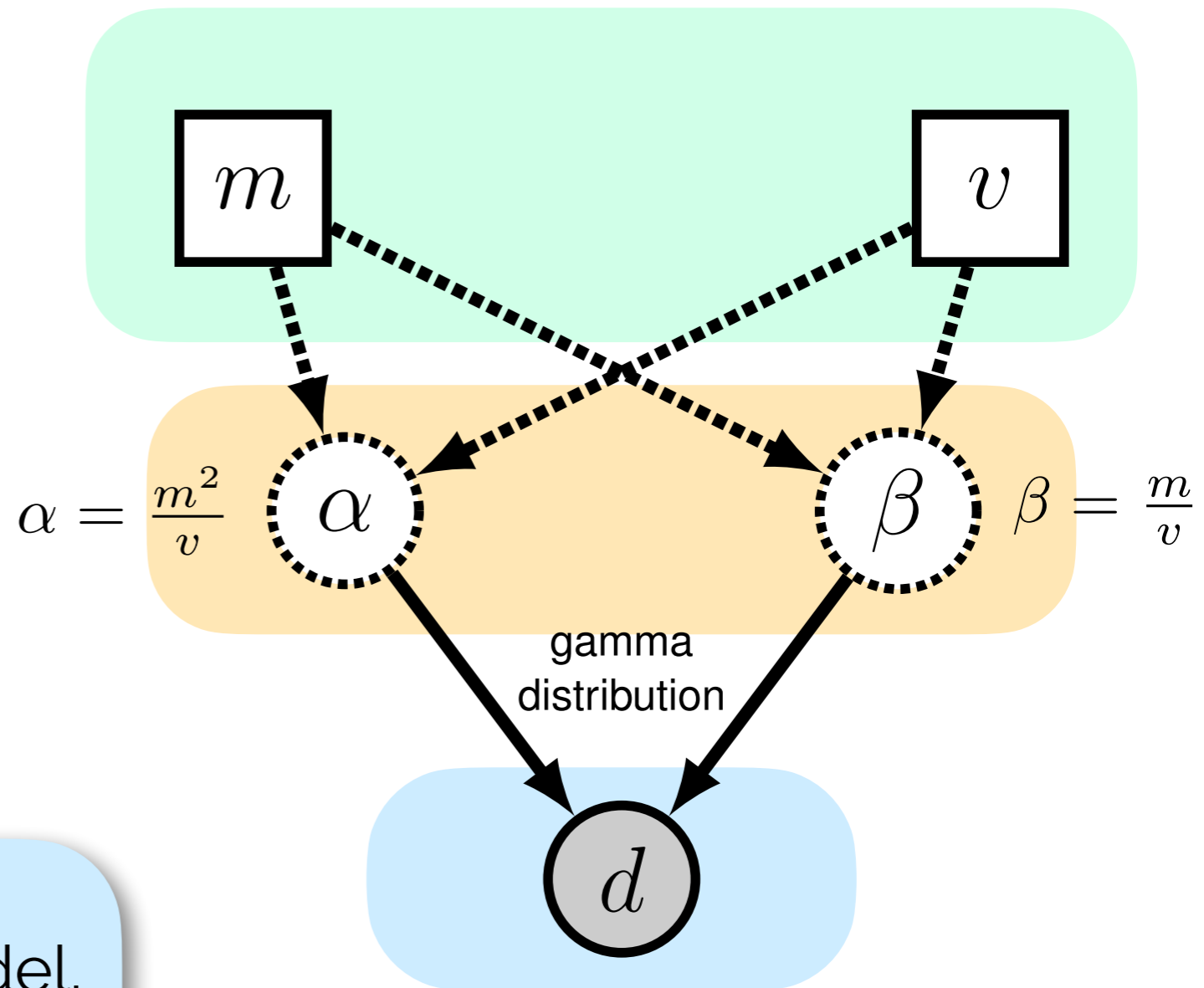
Priors: Archery



Constant nodes represent a fixed value that is asserted or known

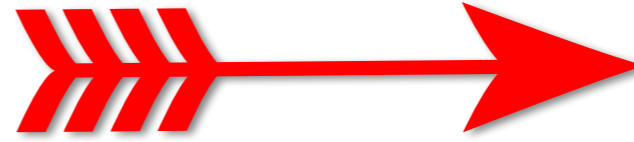
Deterministic nodes represent unknown random variable whose values are determined by other nodes

Stochastic nodes are random variables generated by the model. If we observe the value of a stochastic node, we fix it to that value

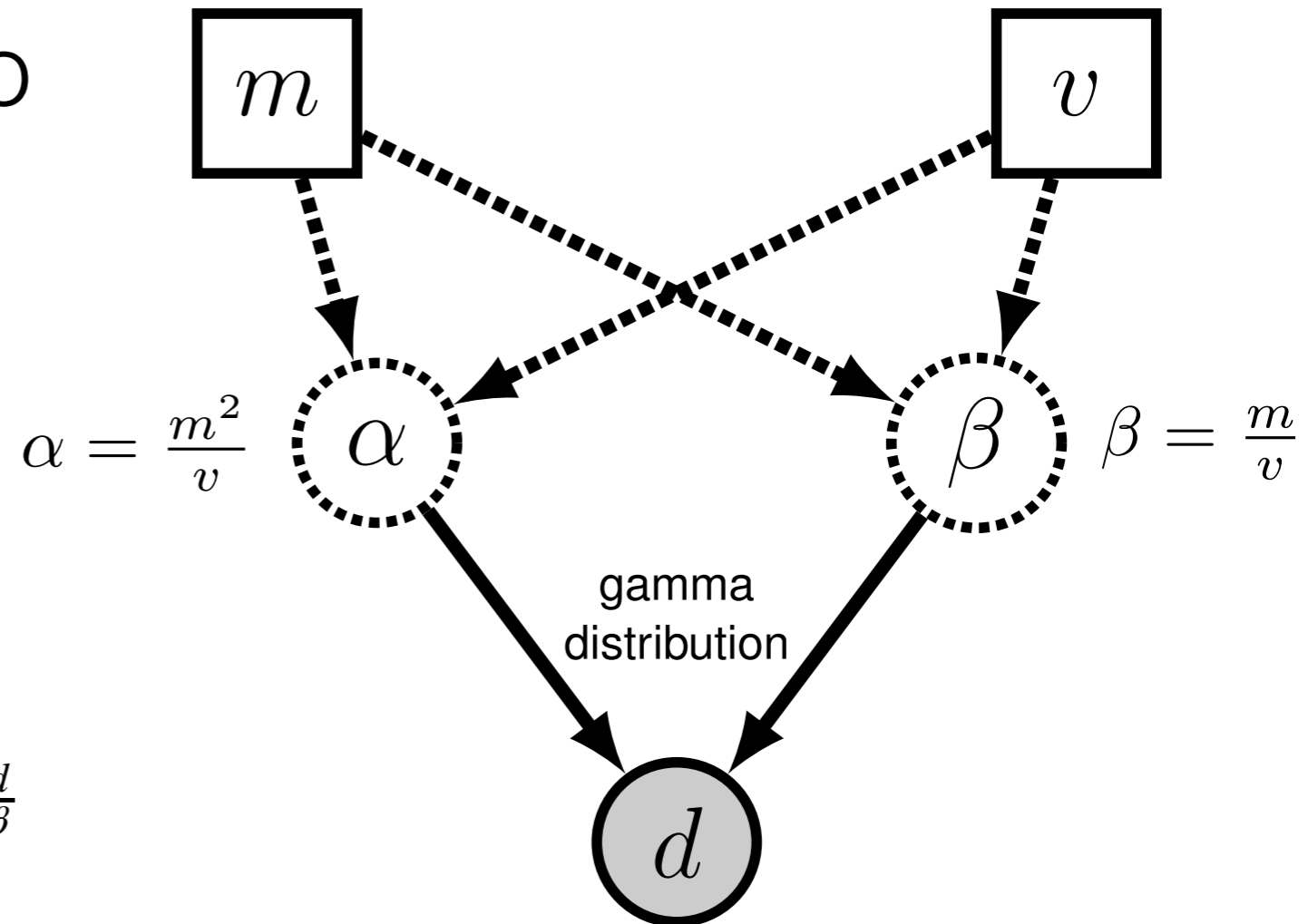


This graphical model has 3 types of nodes

Priors: Archery



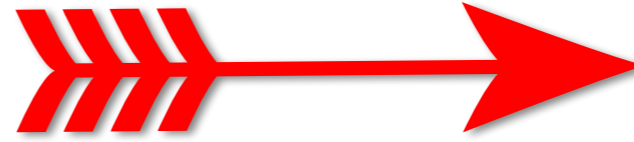
If we set m and v to values corresponding to our assumed model, then we can calculate the likelihood of any observed shot



$$f(d \mid \alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta^\alpha} d^{\alpha-1} e^{-\frac{d}{\beta}}$$

$$f(d = 39.76 \mid \alpha = 1200, \beta = 20) = 7.89916e - 40$$

Priors: Archery



What if we do not know m and v ?

We can use maximum likelihood or Bayesian methods to estimate their values

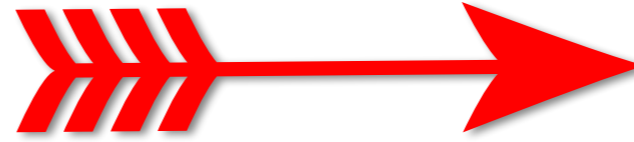
Maximum likelihood methods require us to find the values of m and v that maximize

$$f(d \mid m, v)$$

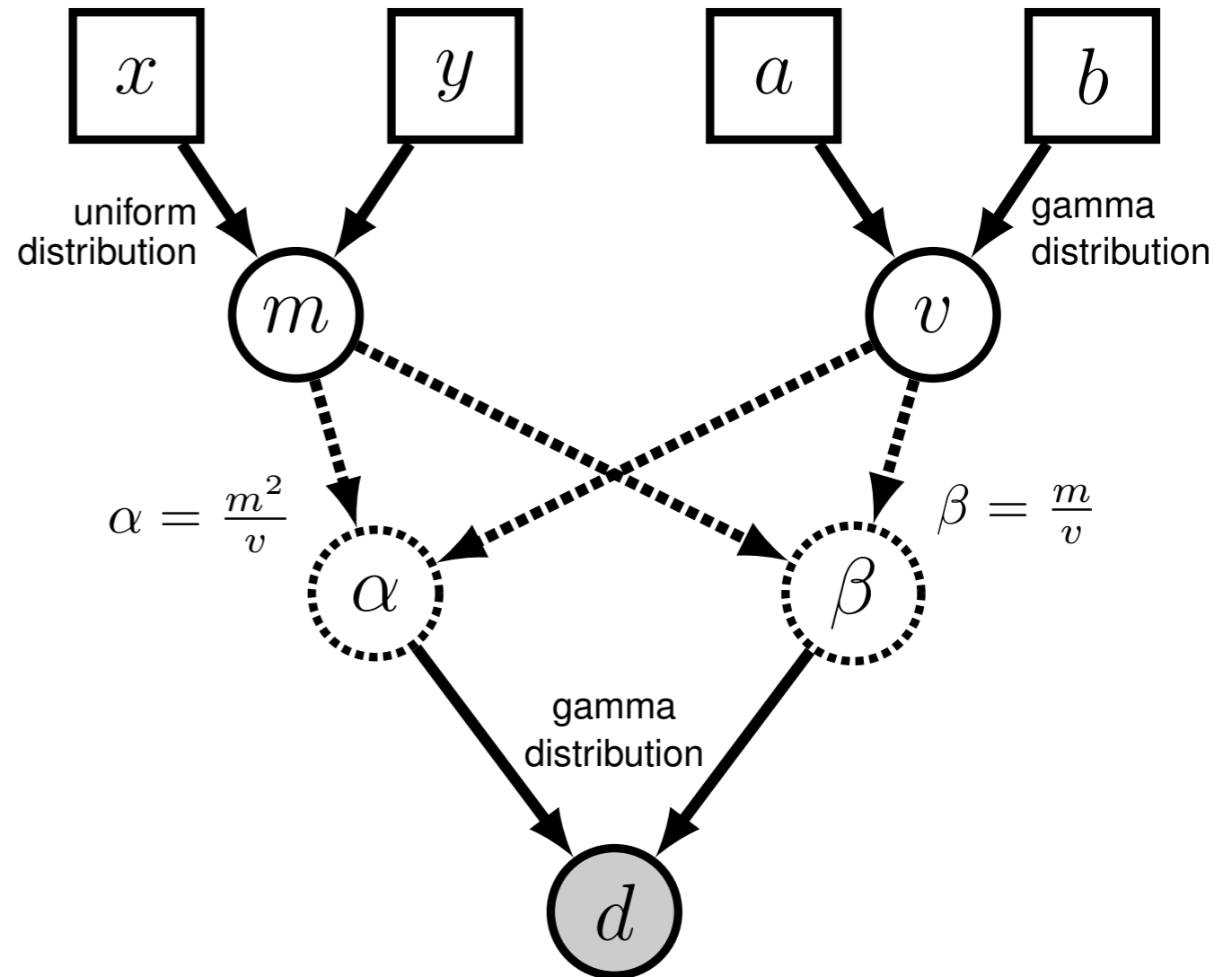
Bayesian methods use prior distributions to describe our uncertainty in m and v and estimate

$$f(m, v \mid d)$$

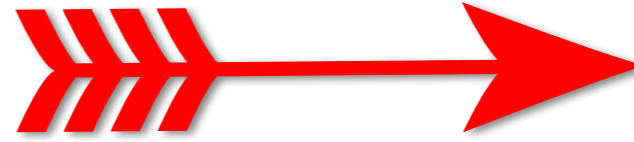
Priors: Archery



We must define prior distributions for m and v to account for uncertainty and estimate the posterior densities of those parameters

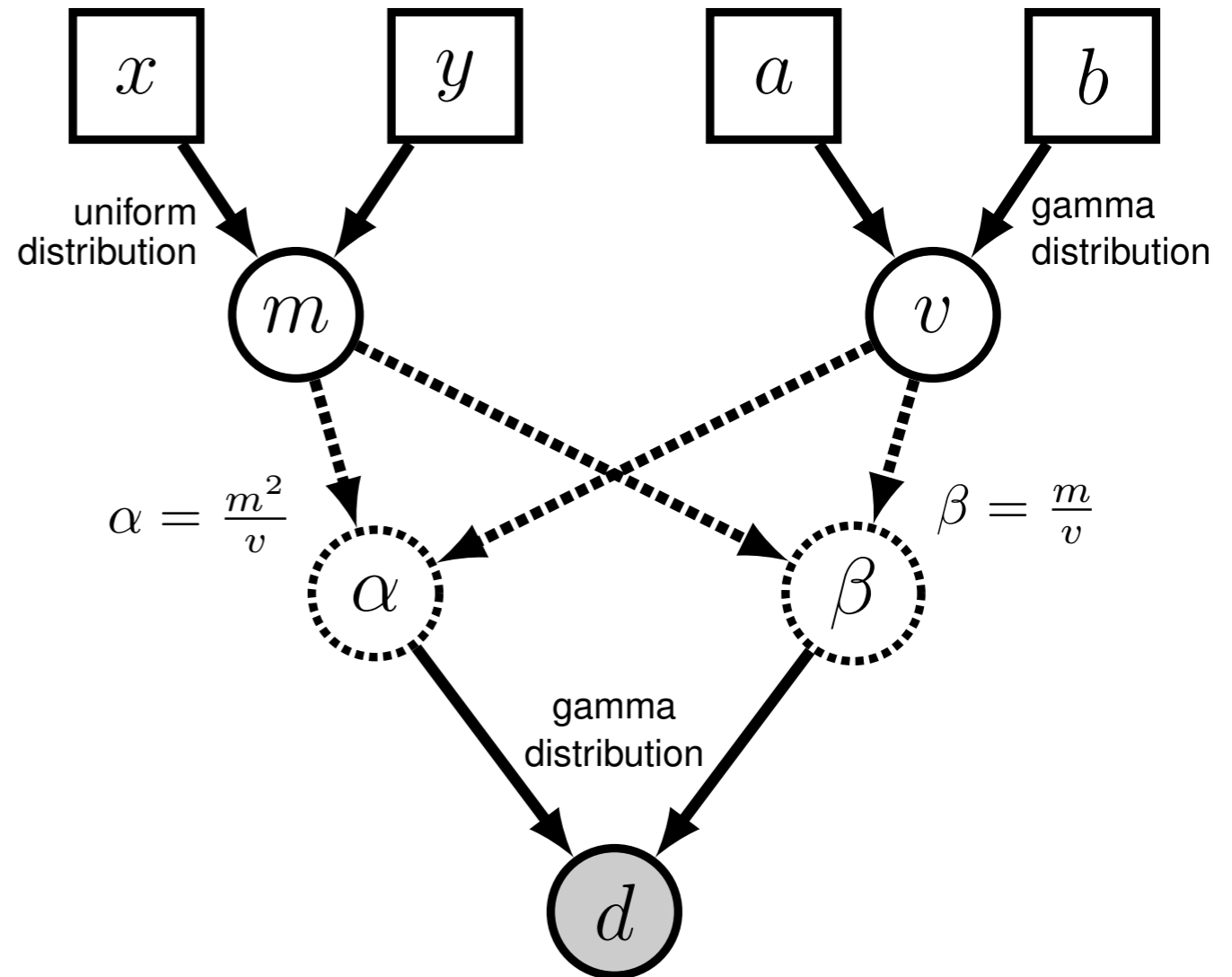


Priors: Archery



Now x and y are the parameters of the uniform prior on m

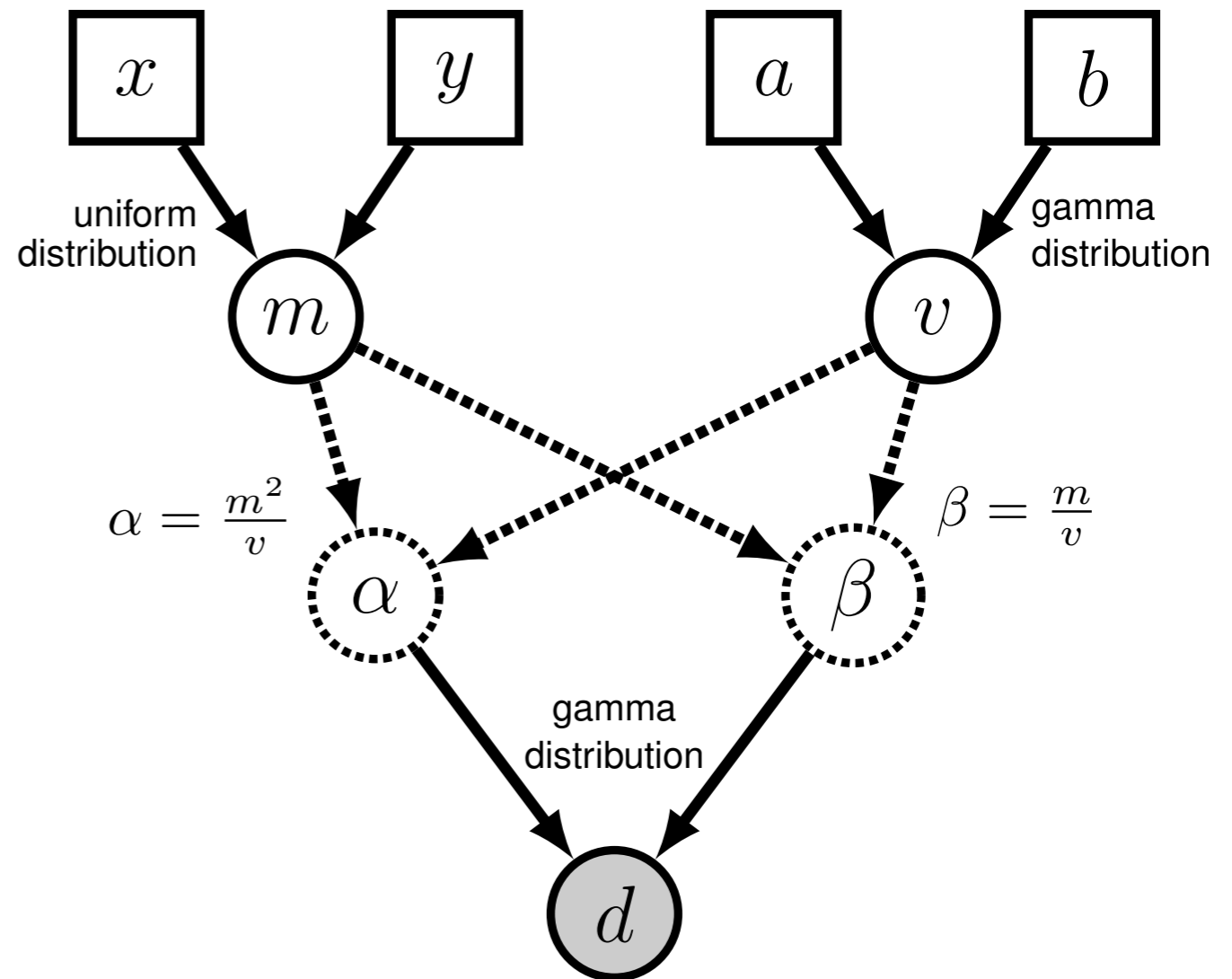
And a and b are the shape and rate parameters of the gamma prior on v



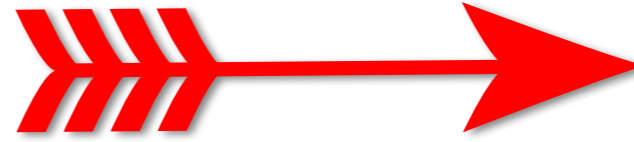
Priors: Archery



Stochastic nodes that are not observed are random variables that are unknown and estimated

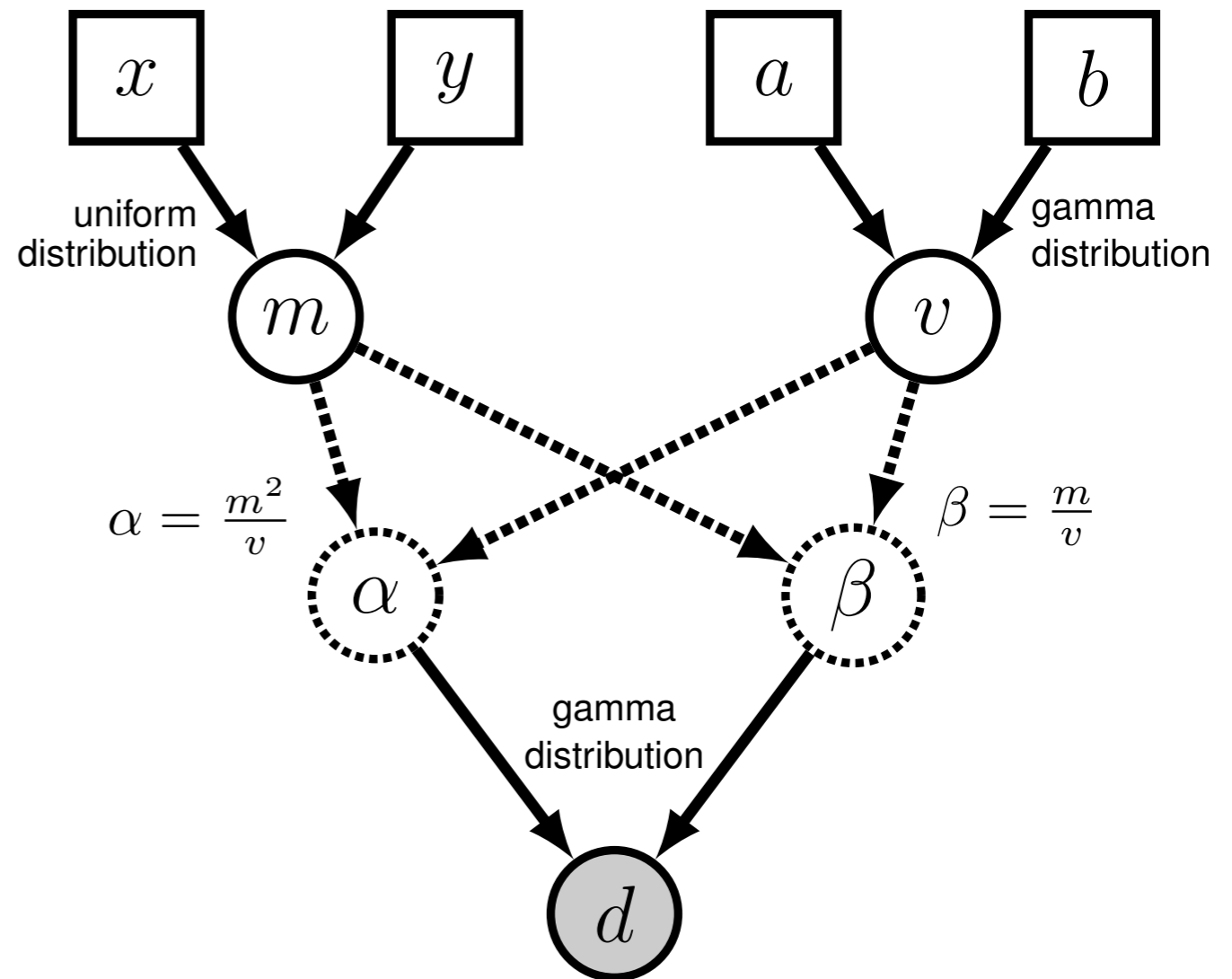


Priors: Archery

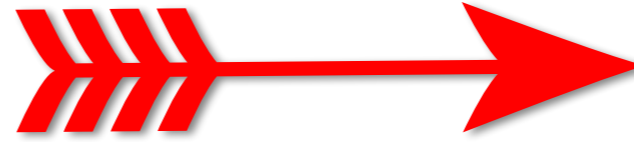


The values we choose for the parameters of these prior distributions should reflect our prior knowledge

If we observed a previous shot at **39.76 cm**, then we can use this to parameterize our priors for analysis of future observations



Priors: Archery



$$m \sim \text{Uniform}(x, y)$$

$$x = 10$$

$$y = 50$$

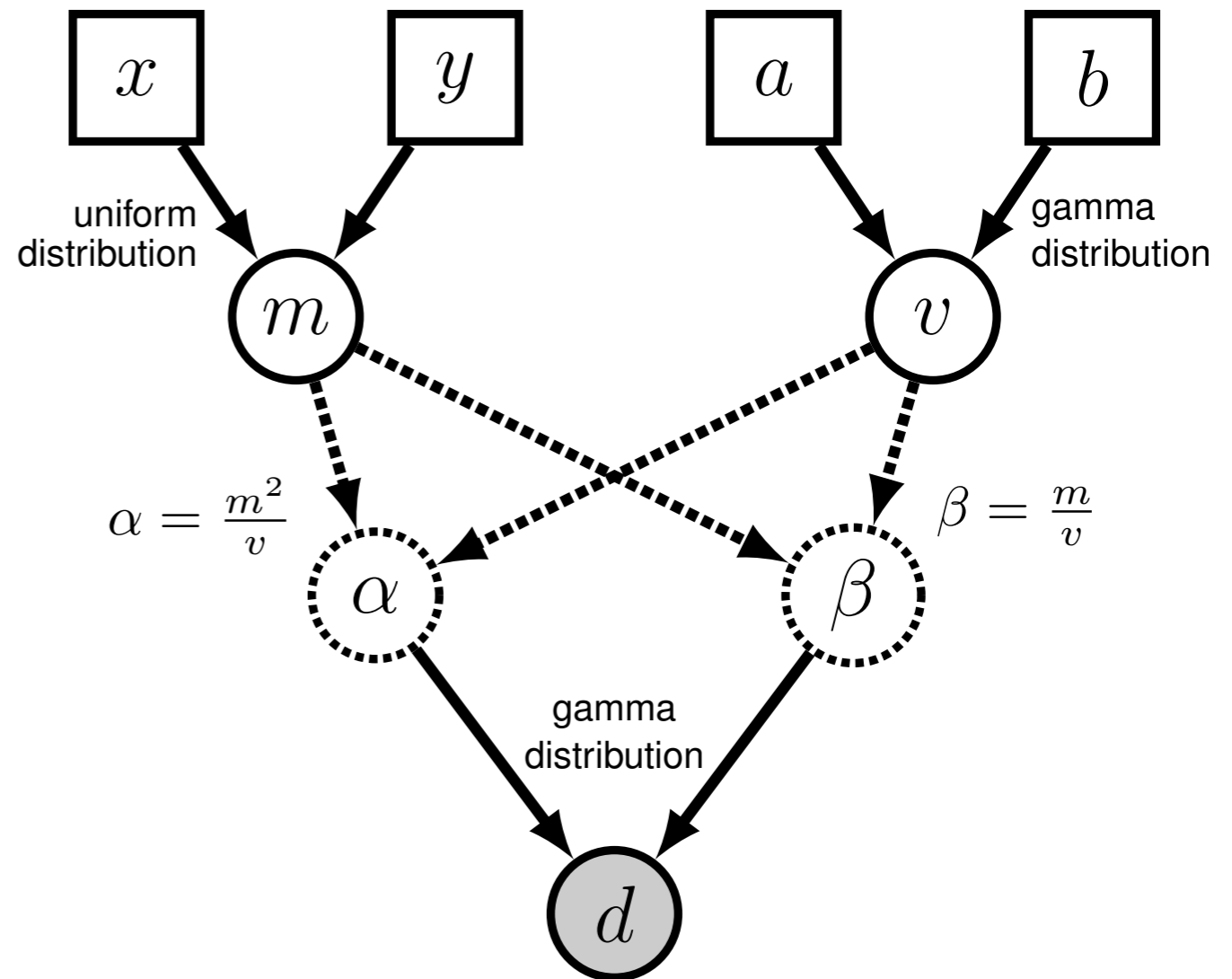
$$\mathbb{E}(m) = 30$$

$$v \sim \text{Gamma}(a, b)$$

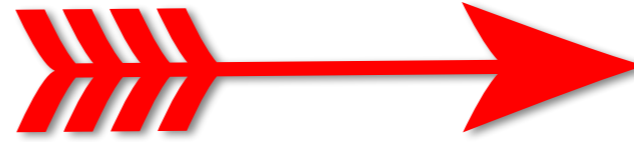
$$a = 20$$

$$b = 2$$

$$\mathbb{E}(v) = 10$$



Priors: Archery

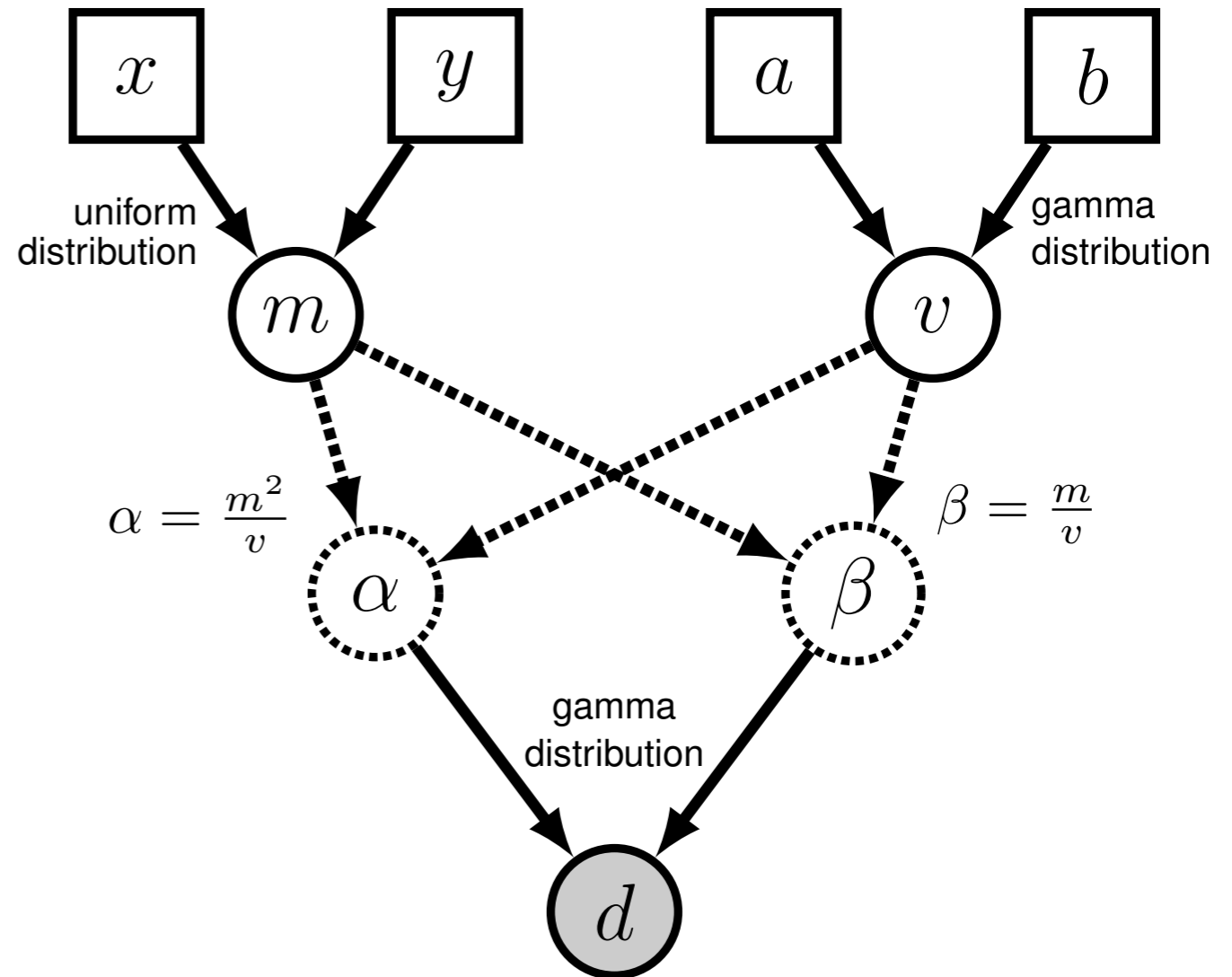


Now that we have a defined model, how do we estimate the posterior probability density?

$$m \sim \text{Uniform}(x, y)$$

$$v \sim \text{Gamma}(a, b)$$

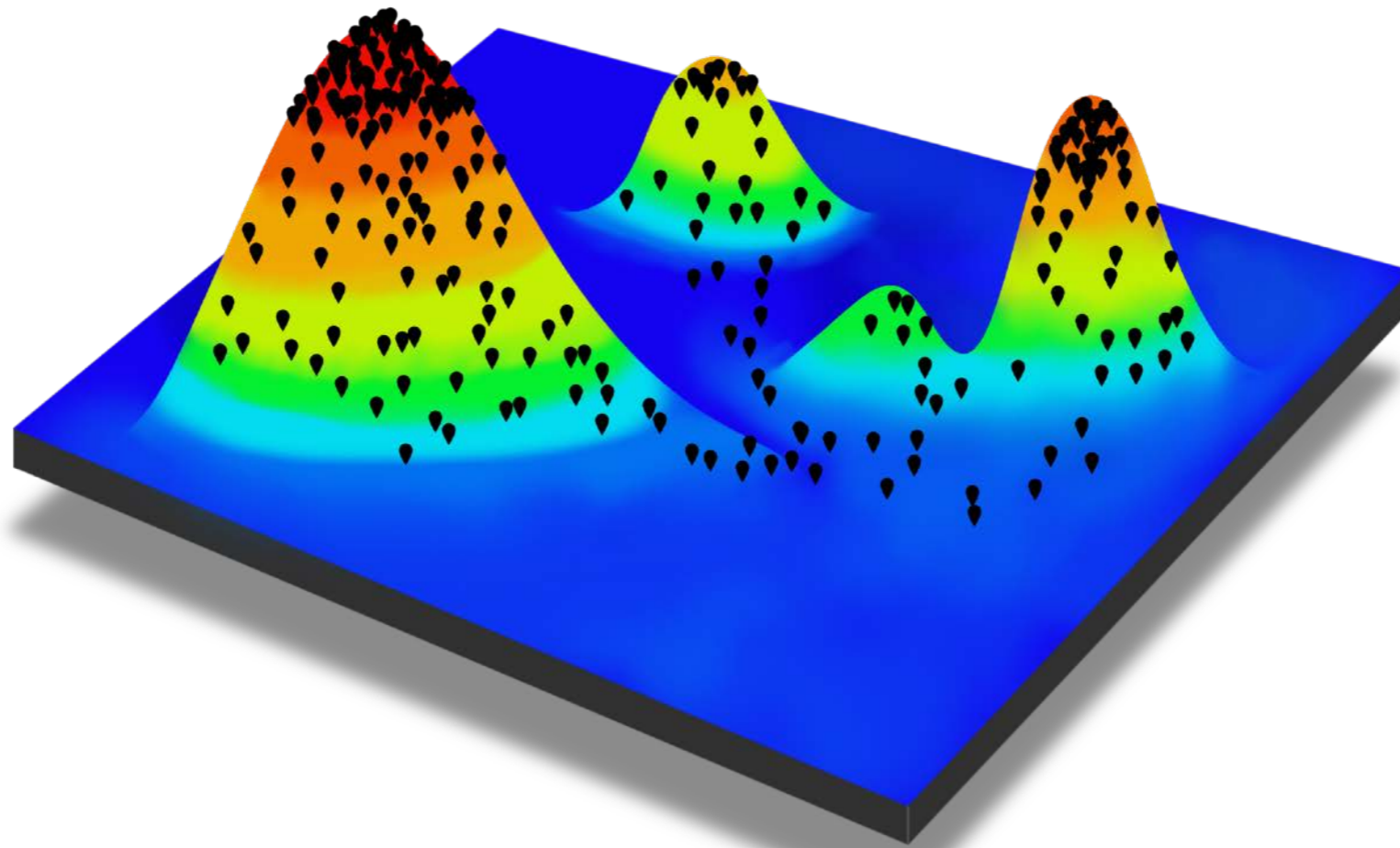
$$d \sim \text{Gamma}(\alpha, \beta)$$



$$f(m, v | d, a, b, x, y) \propto f(d | , \alpha = \frac{m^2}{v}, \beta = \frac{m}{v}) f(m | x, y) f(v | a, b)$$

Markov Chain Monte Carlo

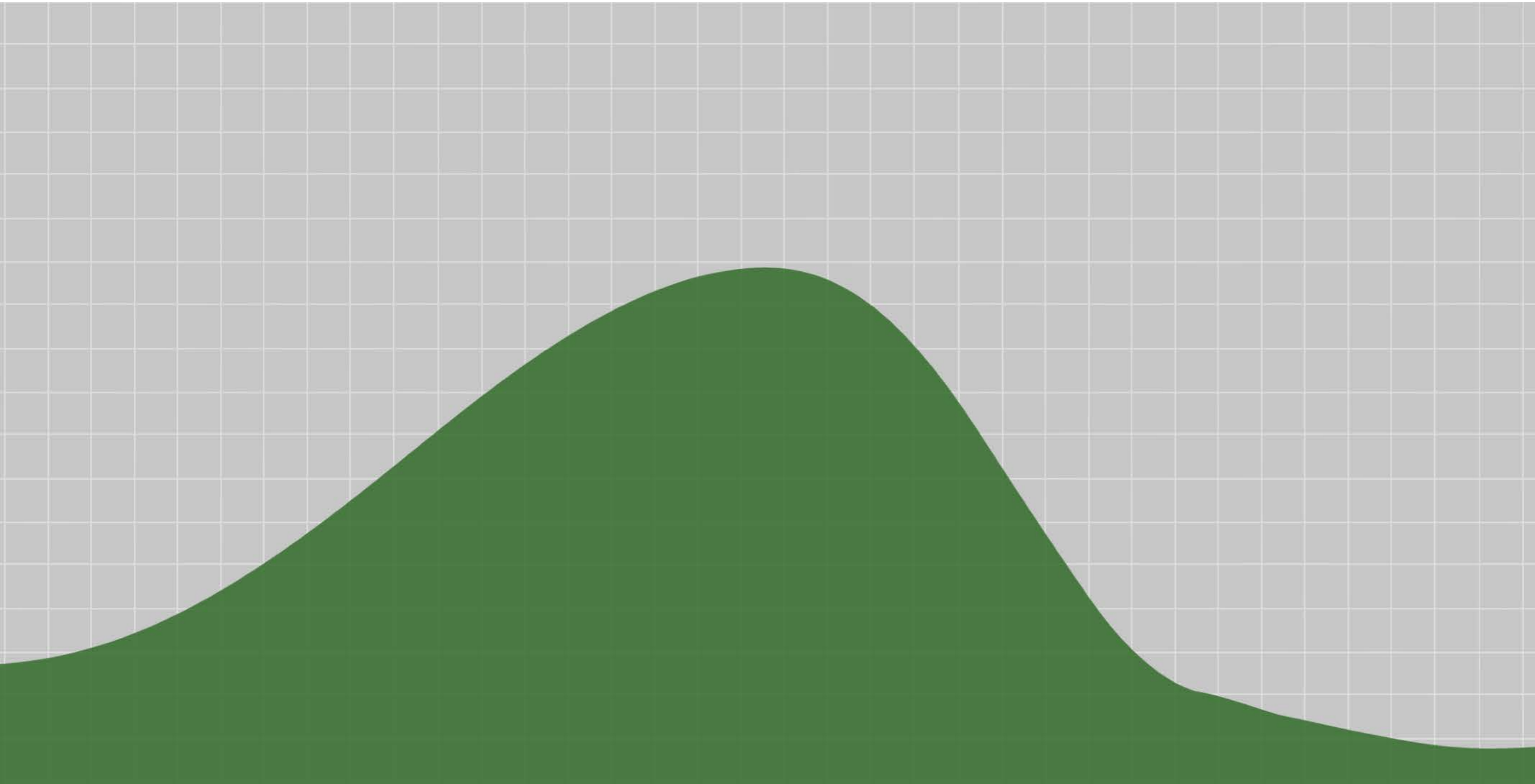
An algorithm for approximating the posterior distribution



Metropolis, et al. 1953. Equations of state calculations by fast computing machines. [J. Chem. Phys.](#)

Hastings. 1970. Monte Carlo sampling methods using Markov chains and their applications. [Biometrika.](#)

MCMC Robot Rules

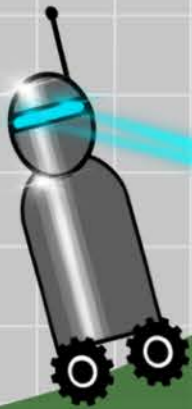


MCMC Robot Rules



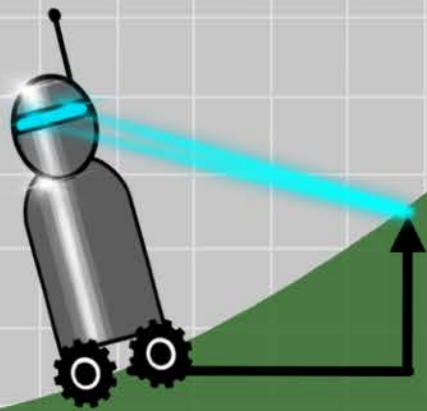
MCMC Robot Rules

uphill moves are
always accepted

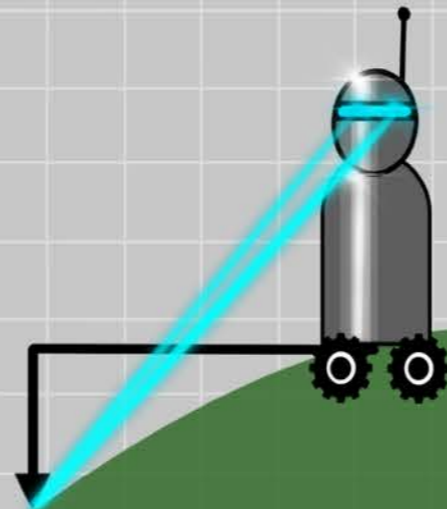


MCMC Robot Rules

uphill moves are
always accepted



MCMC Robot Rules



slightly downhill
moves are usually
accepted

MCMC Robot Rules

extreme downhill
moves are almost
never accepted

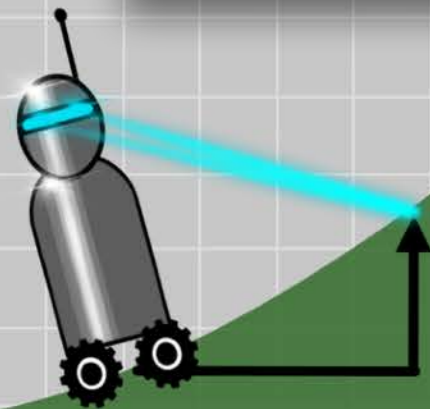


Actual Rules (Metropolis Algorithm)

Metropolis et al. 1953. Equation of state calculations by fast computing machines. *J. Chem. Physics.*

when the robot makes a move
it evaluates the new value
compared to its current state

current = 3.5 m
proposed = 5.5 m
 $R = 5.5 / 3.5 = 2.0$



to decide if it should move to
the proposed state it draws a
variable (x) from a uniform
distribution between 0 and 1

$x \sim \text{Uniform}(0,1)$

the robot moves if:
 $x \leq R$

uphill moves are
always accepted

6
4
2
10
8
6
4
2

Actual Rules (Metropolis Algorithm)

Metropolis et al. 1953. Equation of state calculations by fast computing machines. *J. Chem. Physics.*

current = 11.8 m
proposed = 9.5 m
 $R = 9.5 / 11.8 = 0.81$

$x \sim \text{Uniform}(0, 1)$

the robot moves if:
 $x \leq R$

slightly downhill
moves are usually
accepted

16
14
12
10
8
6
4
2

Actual Rules (Metropolis Algorithm)

Metropolis et al. 1953. Equation of state calculations by fast computing machines. *J. Chem. Physics.*

current = 11 m
proposed = 1 m
 $R = 1.0 / 11.0 = 0.09$

$x \sim \text{Uniform}(0, 1)$

the robot moves if:
 $x \leq R$

extreme downhill
moves are almost
never accepted



16
14
12
10
8
6
4
2

Bayes Rule

posterior probability

likelihood

prior probability

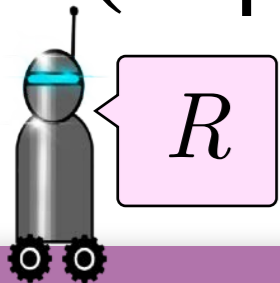
$$\Pr(\theta | D) = \frac{\Pr(D | \theta) \Pr(\theta)}{\sum_{\theta} \Pr(D | \theta) \Pr(\theta)}$$

marginal probability of the data

The diagram illustrates Bayes' Rule. On the left, the posterior probability $\Pr(\theta | D)$ is shown in a yellow box, with an arrow pointing to it from the label "posterior probability". This is followed by an equals sign. To the right of the equals sign is a fraction. The numerator consists of two terms: $\Pr(D | \theta)$ in a blue box, labeled "likelihood" with an arrow, and $\Pr(\theta)$ in a pink box, labeled "prior probability" with an arrow. A horizontal line separates the numerator from the denominator. The denominator is a large green box containing the summation $\sum_{\theta} \Pr(D | \theta) \Pr(\theta)$, which is labeled "marginal probability of the data" with an arrow pointing to it.

Canceling Out the Marginal Likelihood

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



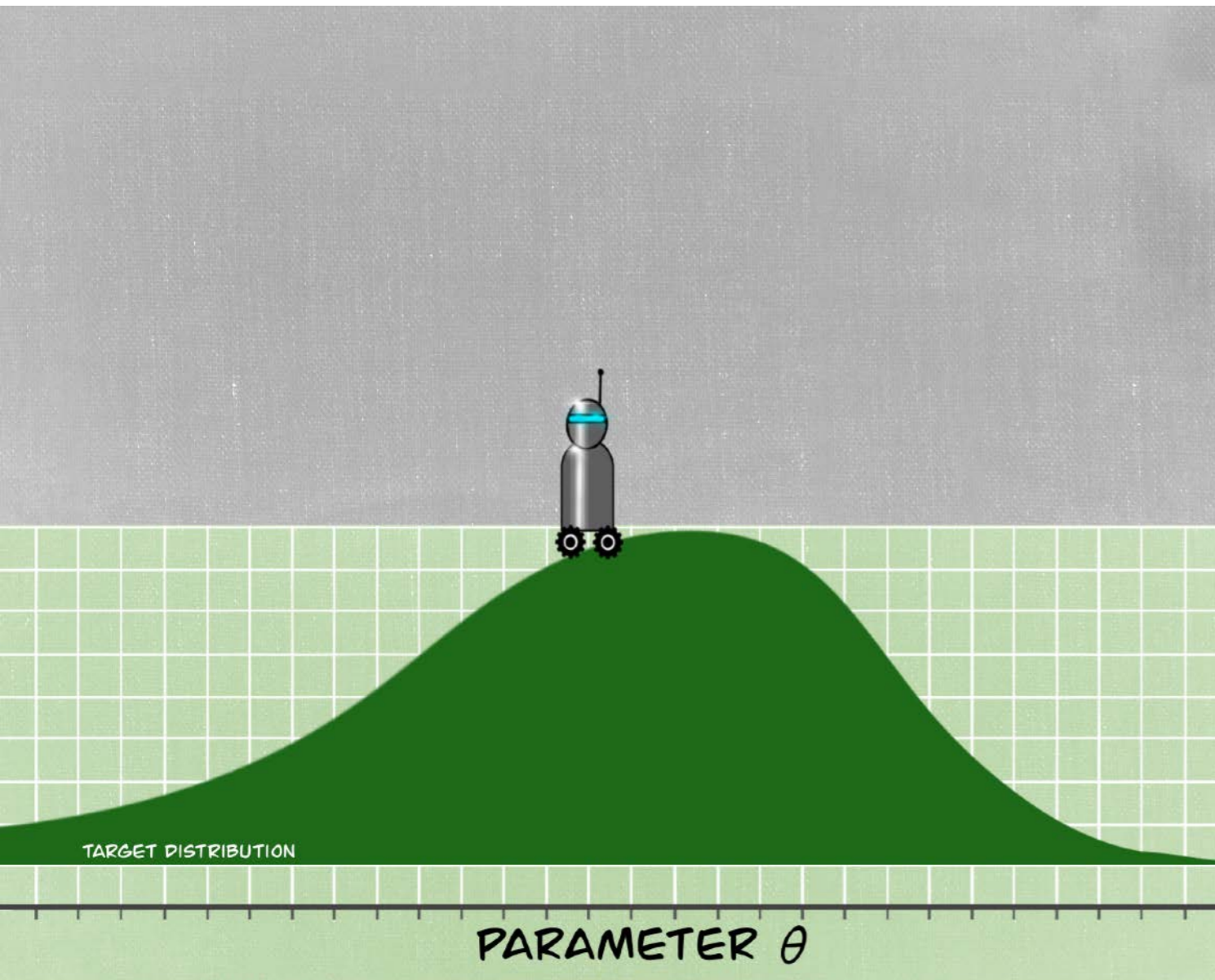
Posterior
odds

Bayes
Rule!

Likelihood
ratio

Prior
odds

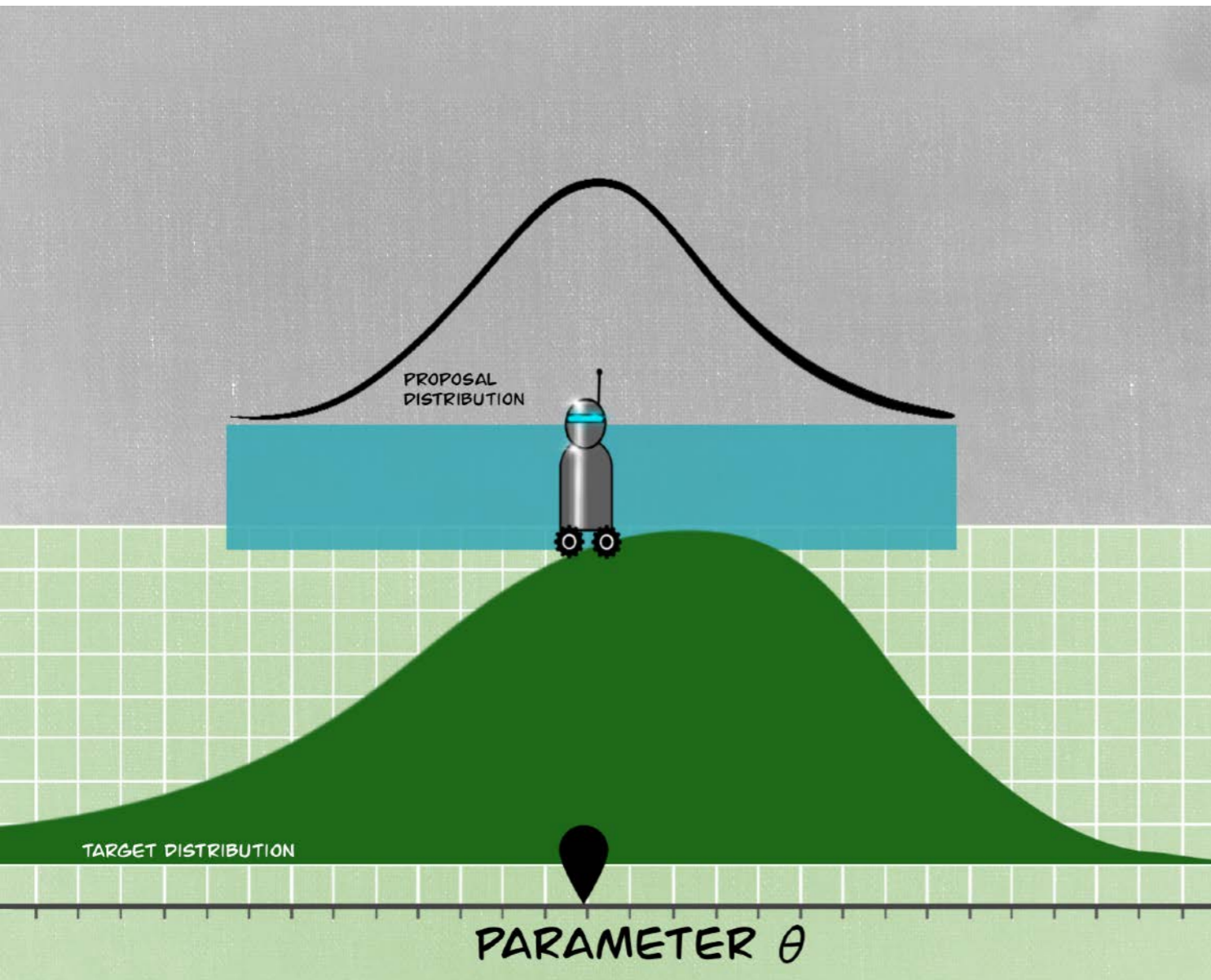
Target and Proposal Distributions



the target distribution is the landscape mapped by the robot

typically, this is the posterior distribution

Target and Proposal Distributions

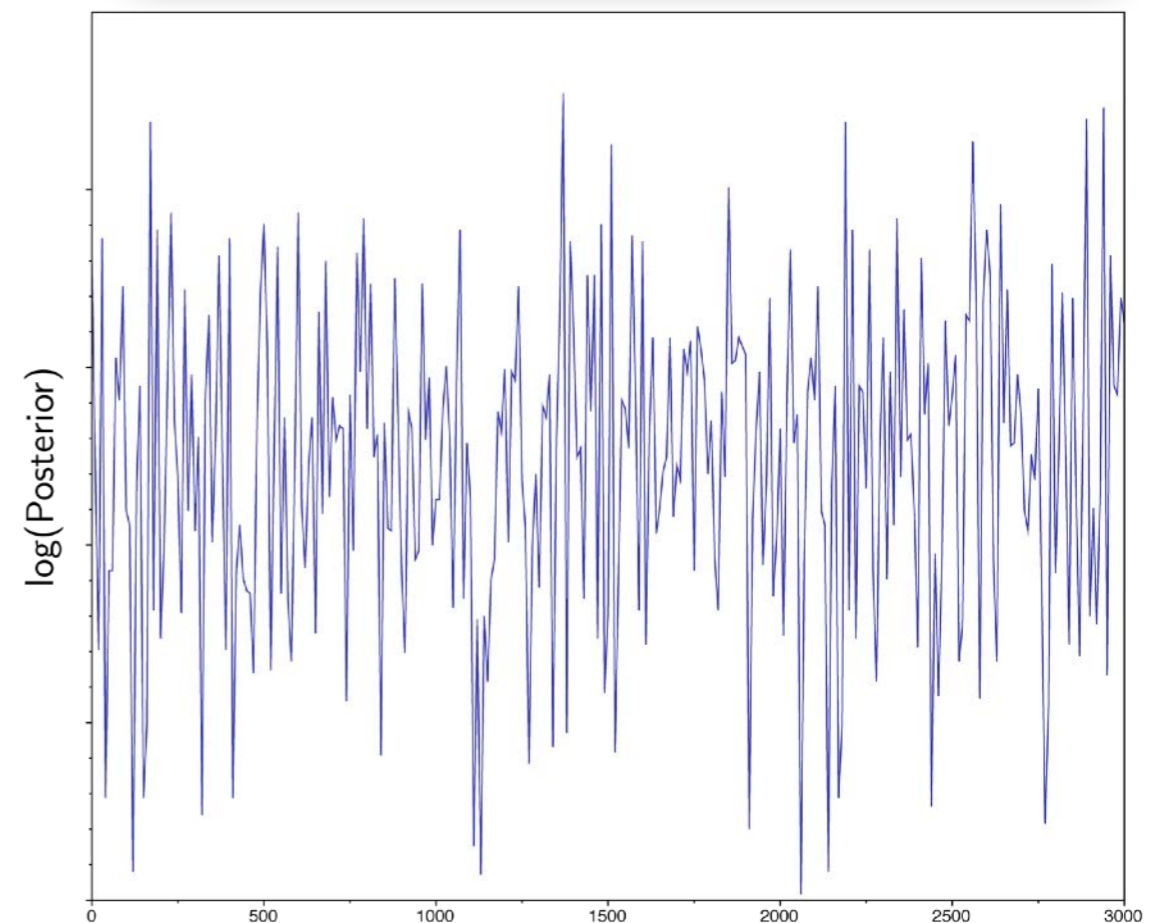
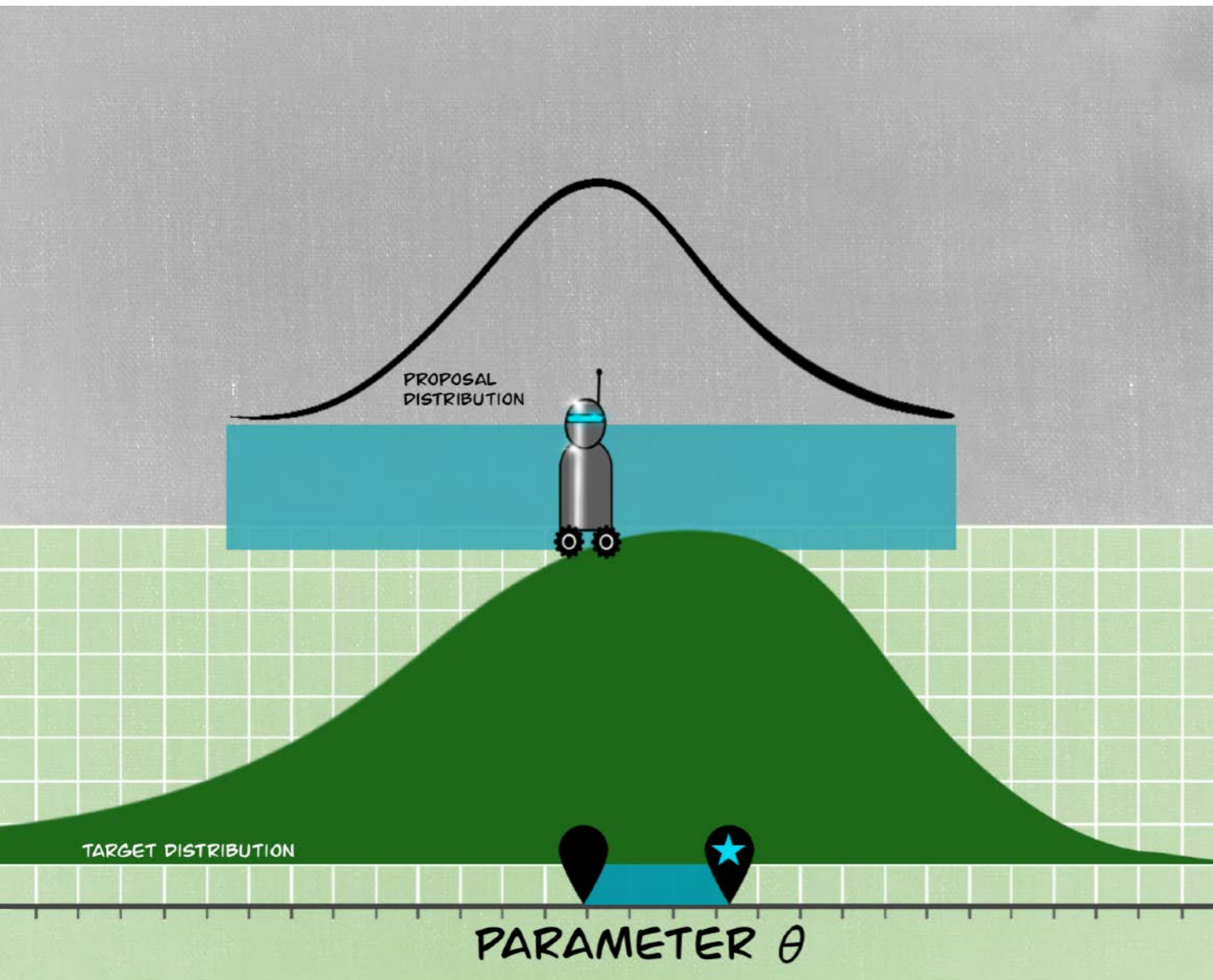


the proposal distribution is separate from the target distribution

the robot uses the proposal distribution to choose the next spot to move

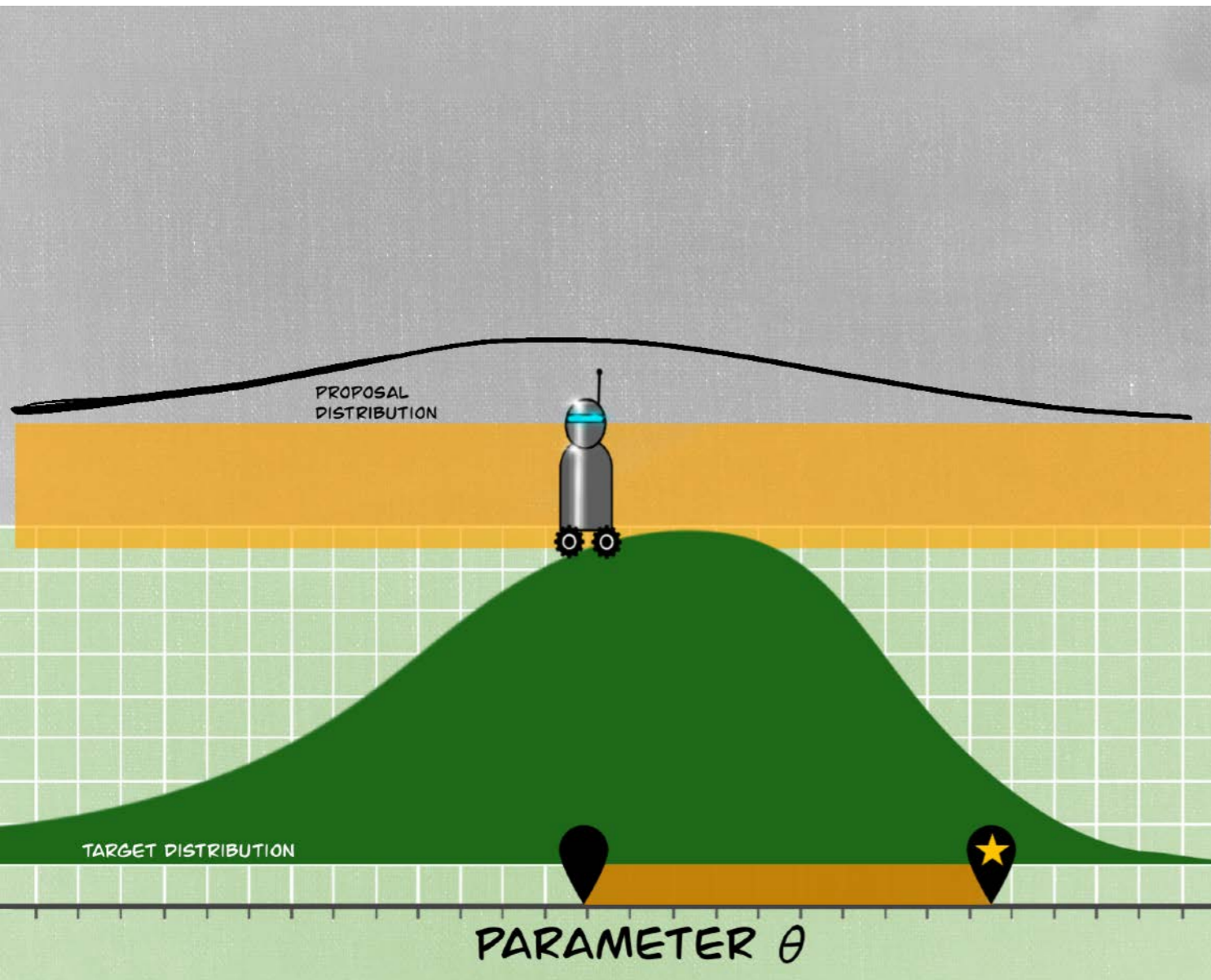
Target and Proposal Distributions

a good proposal distribution samples the target distribution effectively (i.e., "good mixing")

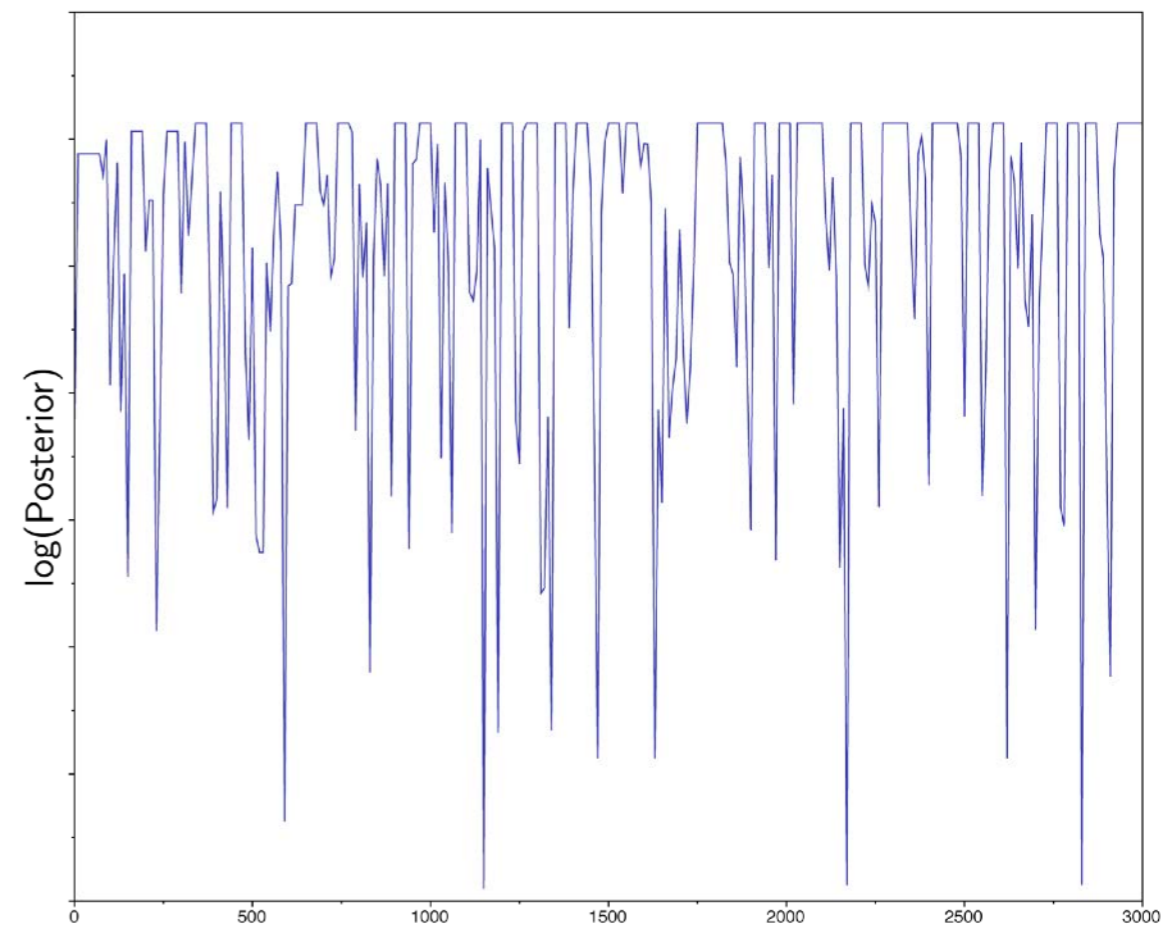


a trace plot of the sampled parameter values looks like white noise

Target and Proposal Distributions



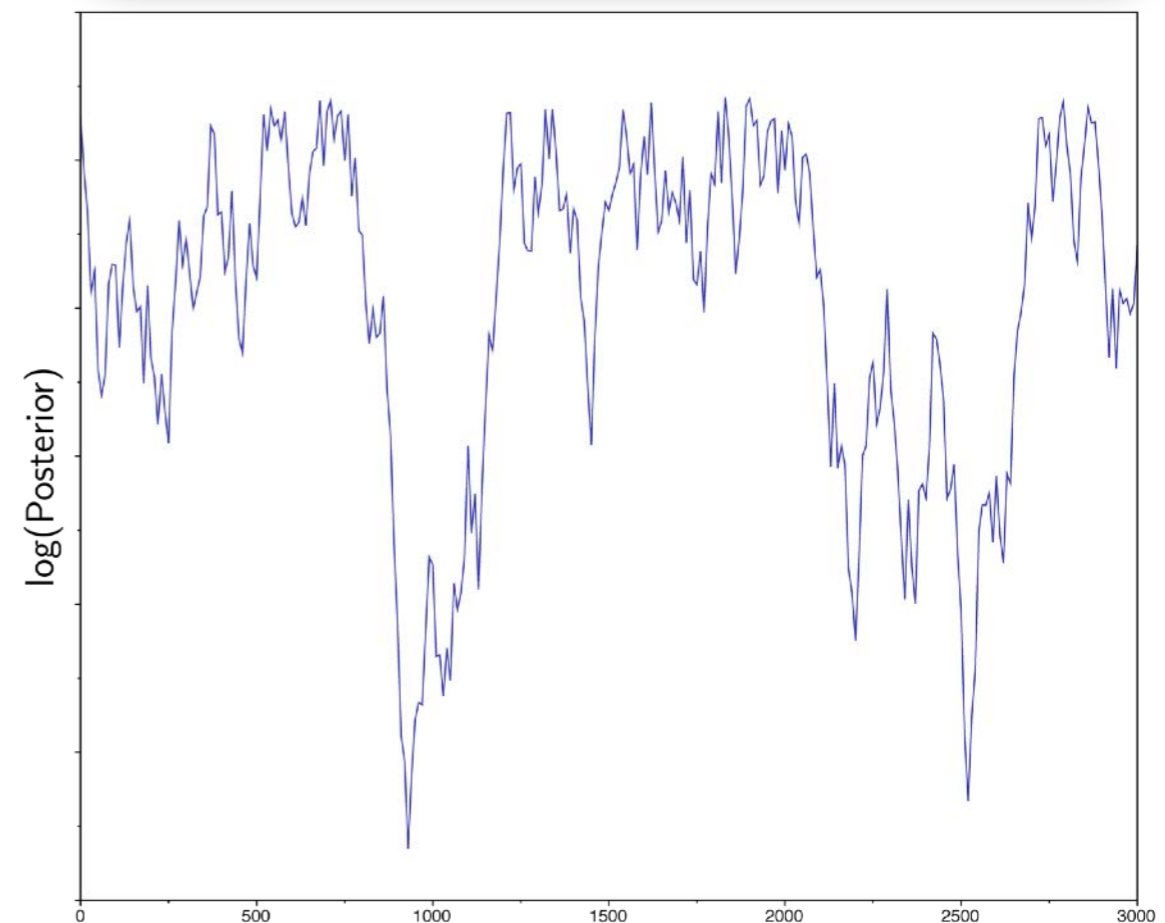
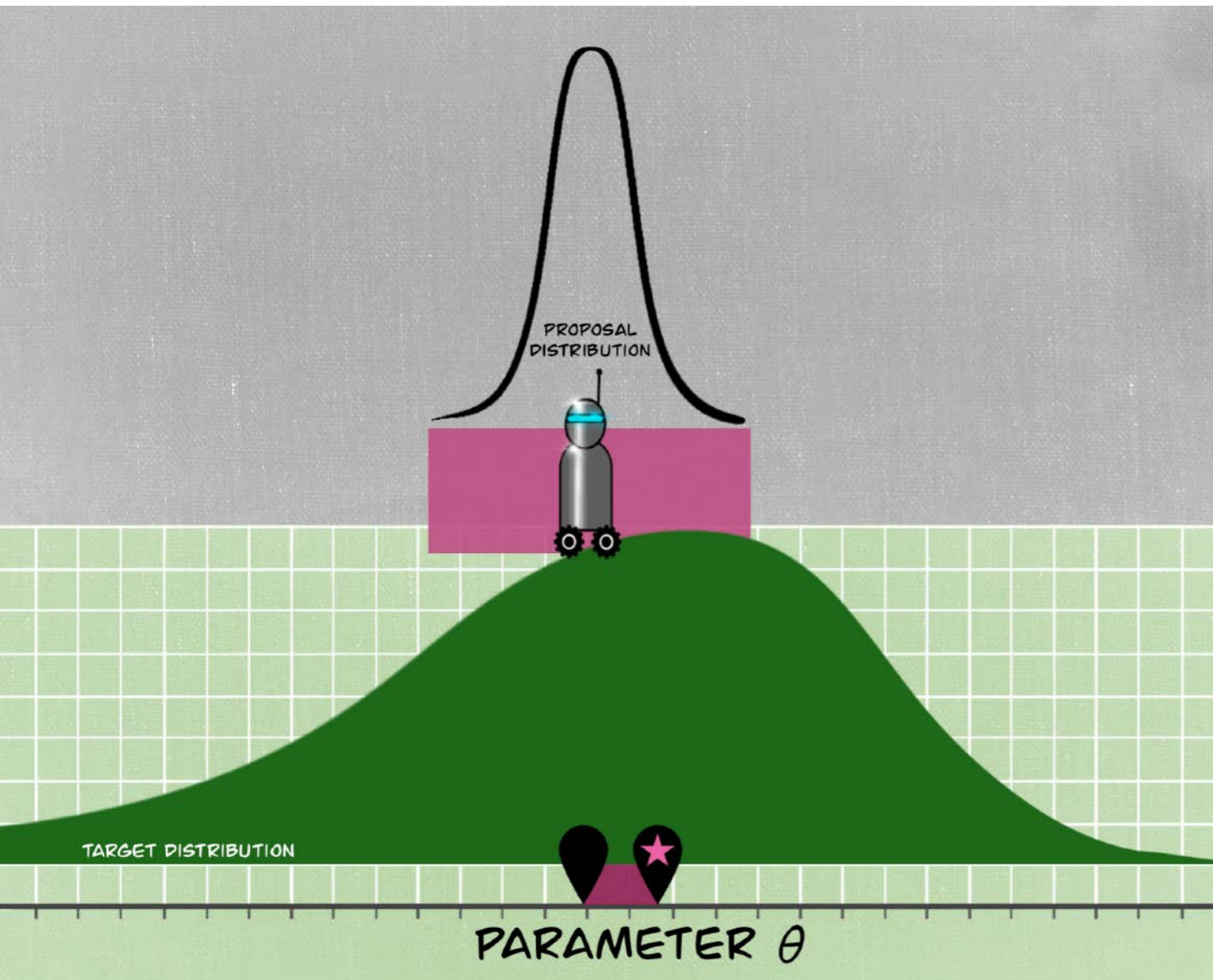
an overly bold proposal results in many rejected moves



this causes the robot to get stuck, seen as plateaus in the trace plot

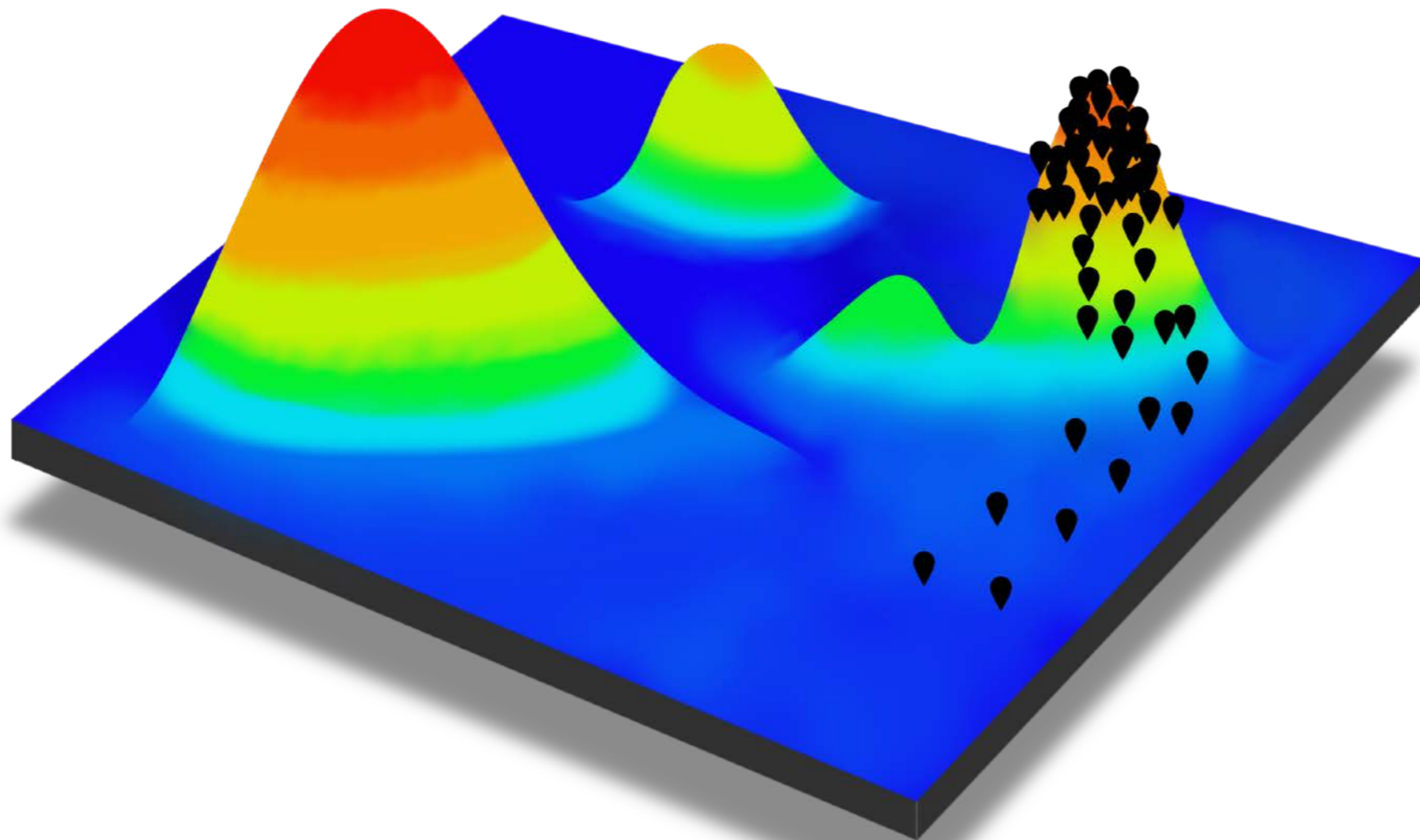
Target and Proposal Distributions

a proposal distribution that only allows for baby steps results in lots of accepted moves



this causes big waves in the trace plot as the robot takes small incremental samples

Metropolis Coupled MCMC



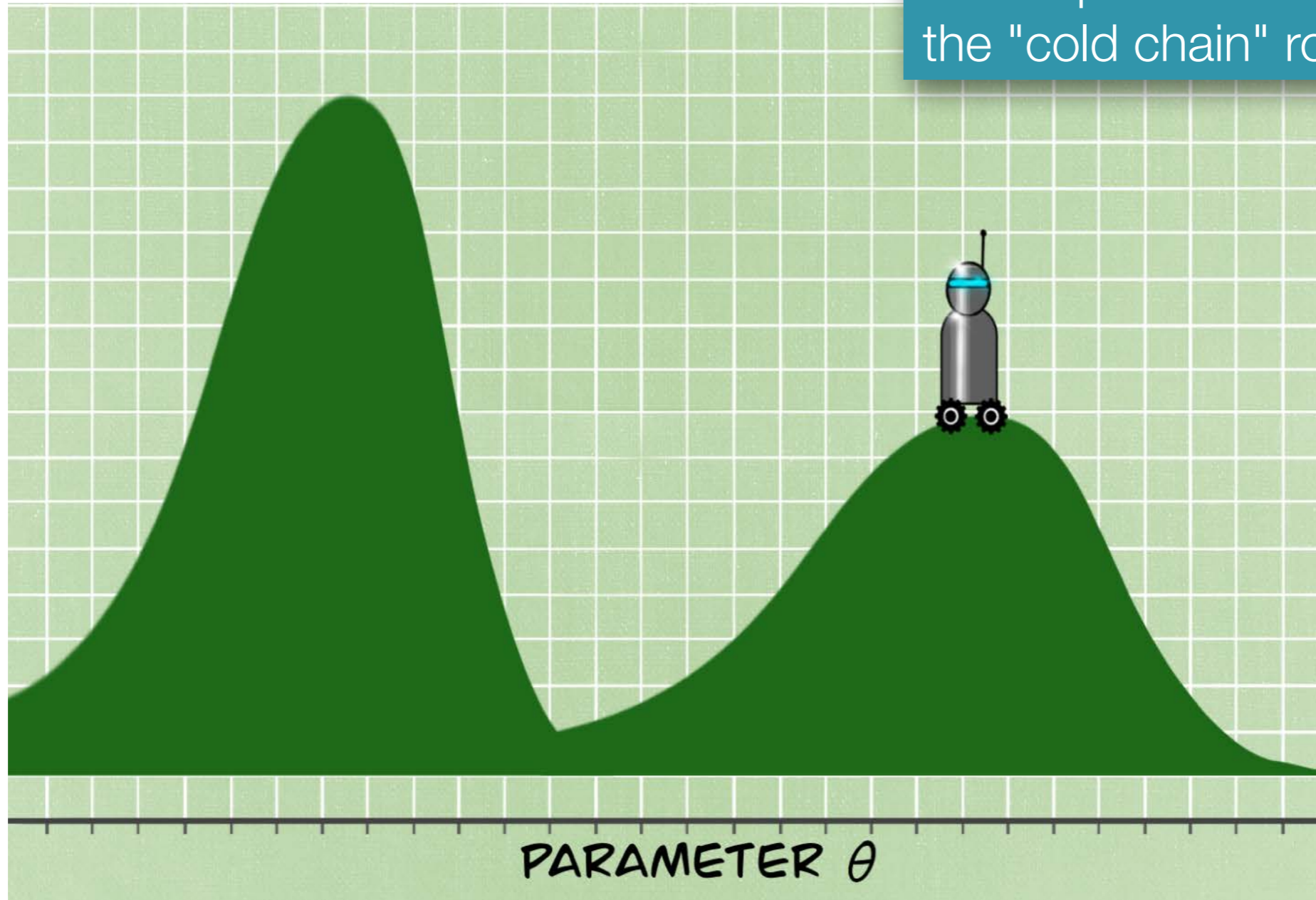
sometimes even
good robots need
help

MCMCMC
introduces helper
robots that act as
scouts to explore
more parameter
space

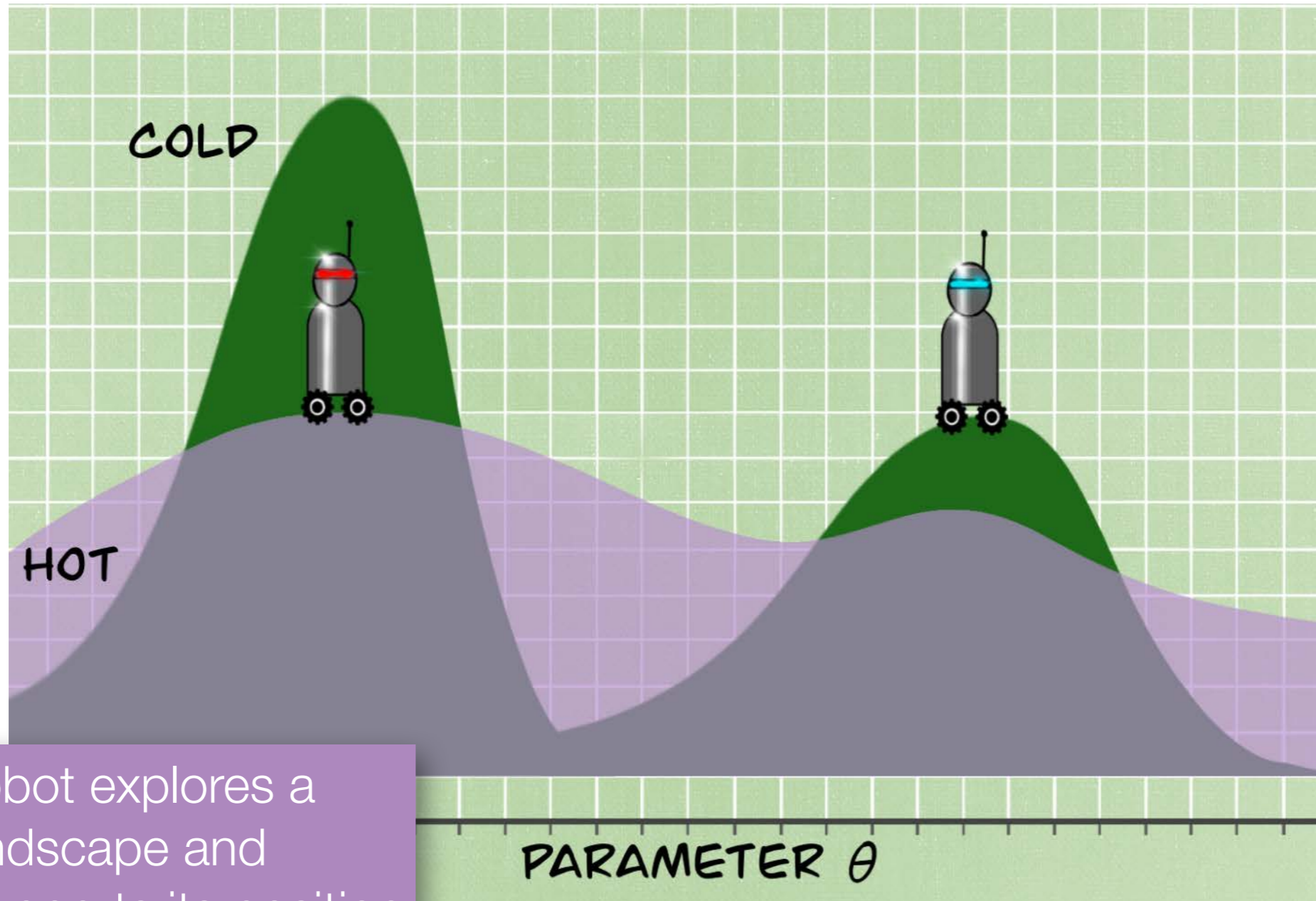
Geyer, C. J. 1991. Markov chain Monte Carlo maximum likelihood for dependent data. Pages 156-163 in Computing Science and Statistics (E. Keramidas, ed.).

Metropolis Coupled MCMC

all samples are taken by the "cold chain" robot

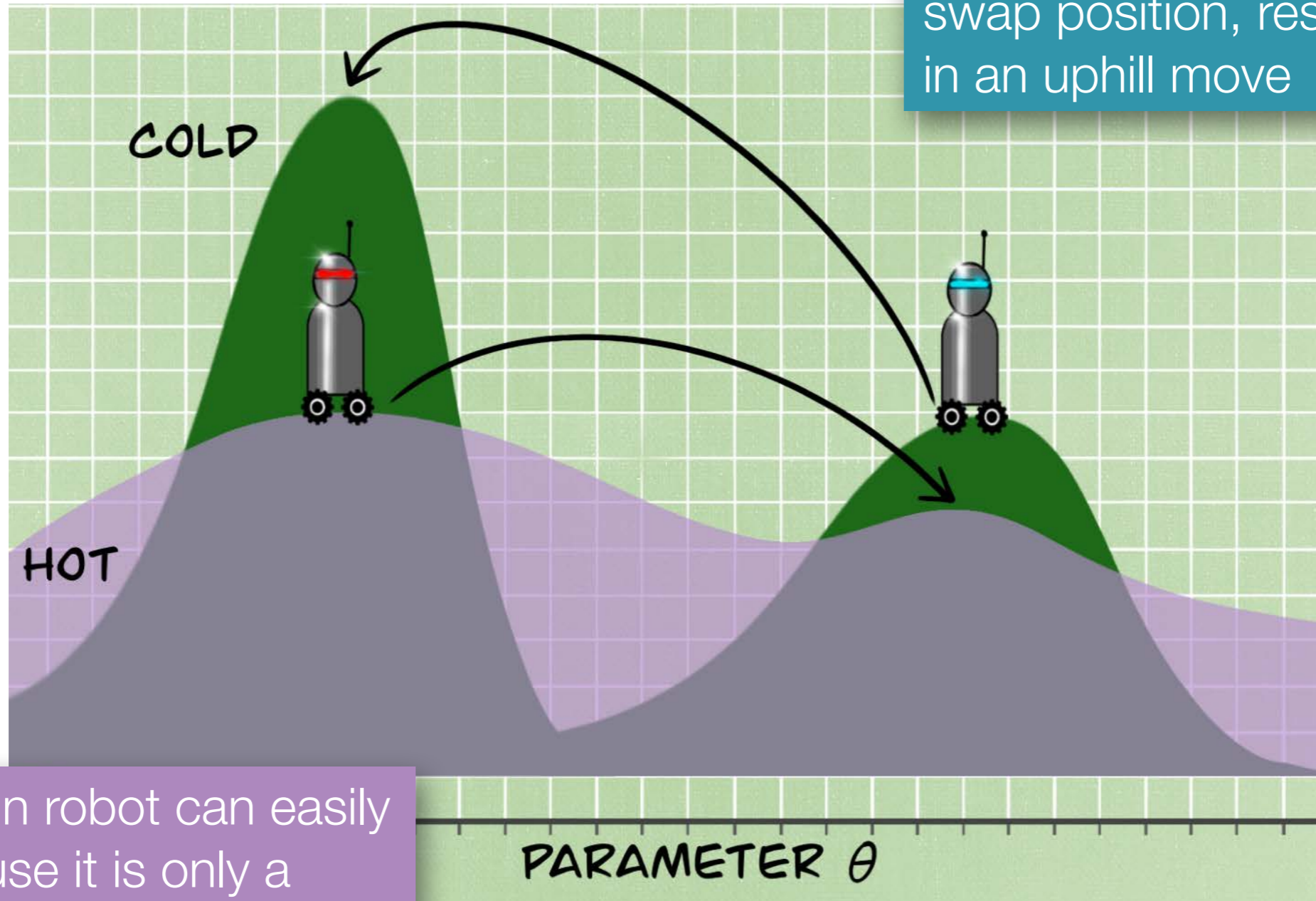


Metropolis Coupled MCMC



the scout robot explores a "heated" landscape and periodically reports its position to the cold chain robot

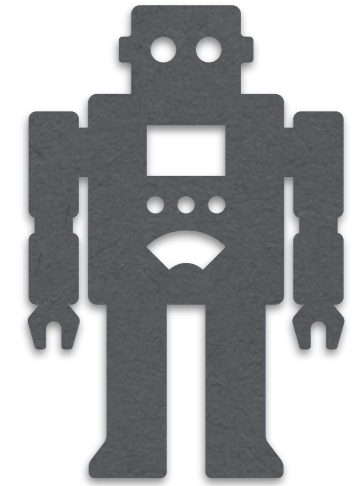
Metropolis Coupled MCMC



the cold-chain robot will swap position, resulting in an uphill move

the hot-chain robot can easily swap because it is only a slightly downhill move

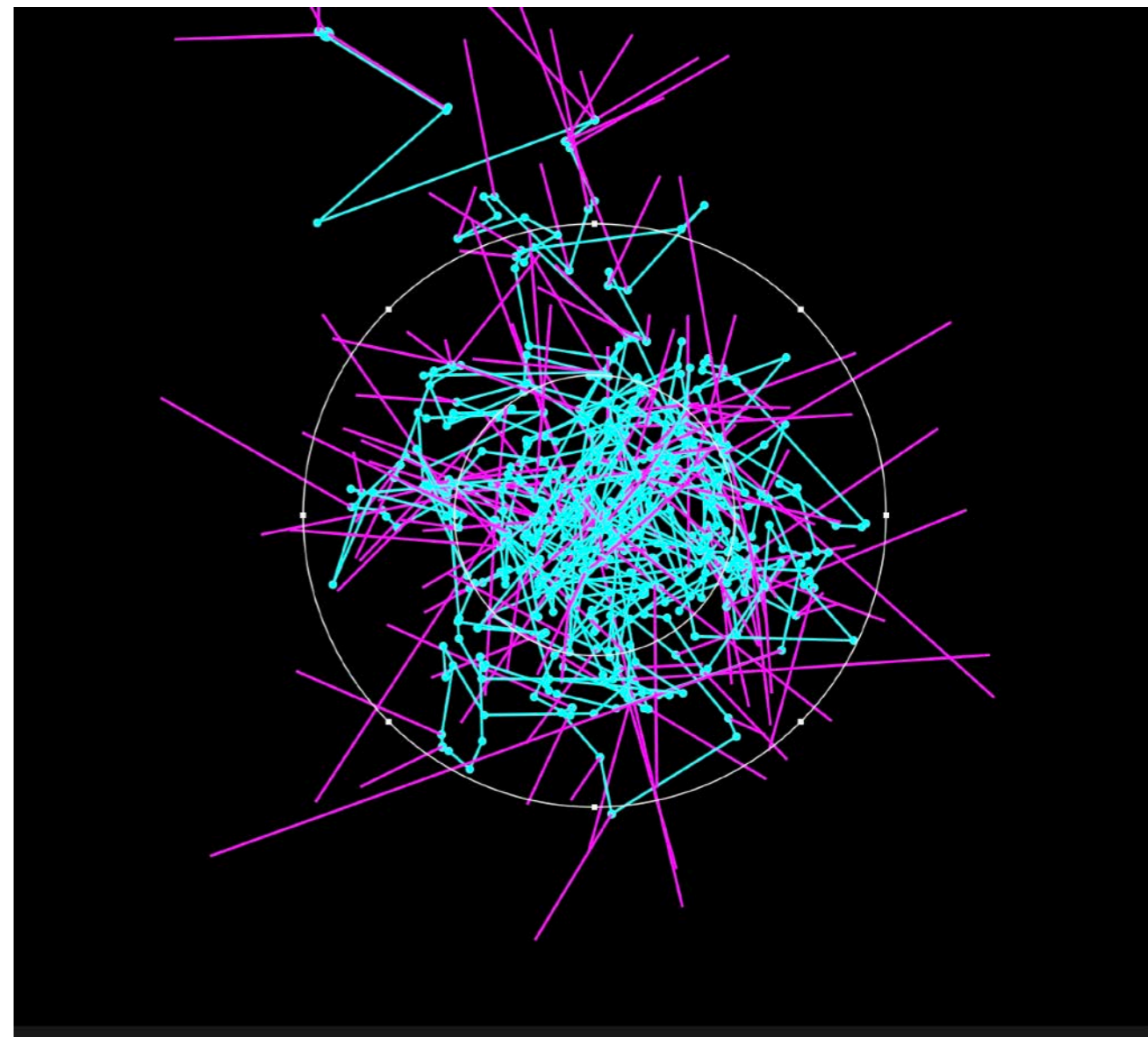
Markov Chain Monte Carlo



Learn more about MCMC!

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

MCMCRobot, a helpful tool for learning MCMC by Paul Lewis



Markov Chain Monte Carlo

Learn more about MCMC!

REVIEW ARTICLE

DOI: 10.1038/s41559-017-0280-x

nature
ecology & evolution

A biologist's guide to Bayesian phylogenetic analysis

Fabrícia F. Nascimento ^{1,4*}, Mario dos Reis ² and Ziheng Yang ^{3*}

<https://thednainus.wordpress.com/2017/03/03/tutorial-bayesian-mcmc-phylogenetics-using-r/>