

INTRODUCTION TO BAYESIAN PHYLOGENETICS USING RevBAYES

Tracy A. Heath & Walker Pett

**Ecology, Evolution, & Organismal Biology
Iowa State University**

Biogeography Meeting
Bengaluru, India
September 25, 2017

YOUR INSTRUCTORS



Tracy Heath



Walker "Will" Pett

WORKSHOP OBJECTIVES

The goal of this workshop is to give you an introduction to Bayesian phylogenetic inference in RevBayes

After today, we hope that you will become familiar with...

- the concepts of Bayesian inference, hierarchical models, and MCMC.
- the syntax of Rev, the model specification language used in RevBayes.
- divergence-time estimation and using fossil data.
- specifying the model and MCMC for a phylogenetic analysis in RevBayes.

WORKSHOP OBJECTIVES

The concepts of Bayesian phylogenetic inference and RevBayes are complex and may be challenging.

If you're planning on using RevBayes for your own analyses, you may still find this difficult after the workshop.

But, you will know...

- where to find documentation (the [tutorials page](#)).
- where to ask for help (the [user forum](#)).
- how to get started with an analysis.
- the many possible types of analyses in RevBayes.

WORKSHOP MATERIALS

Links to all of the workshop materials are available on our website:

<http://phyloworks.org/revbayes-workshop2017>

BAYESIAN OR MAXIMUM LIKELIHOOD?

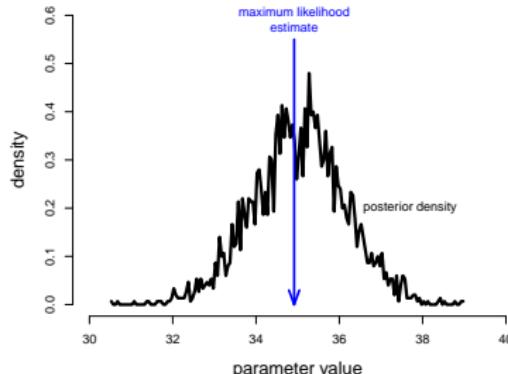
What's the difference?

Bayesian:

- estimates $f(\theta | \mathcal{D})$
- estimates a distribution
- parameters are random variables
- average over nuisance parameters

Maximum Likelihood:

- maximizes $f(\mathcal{D} | \theta)$
- point estimate
- parameters are fixed/unknown
- optimize nuisance parameters



JOINT PROBABILITIES

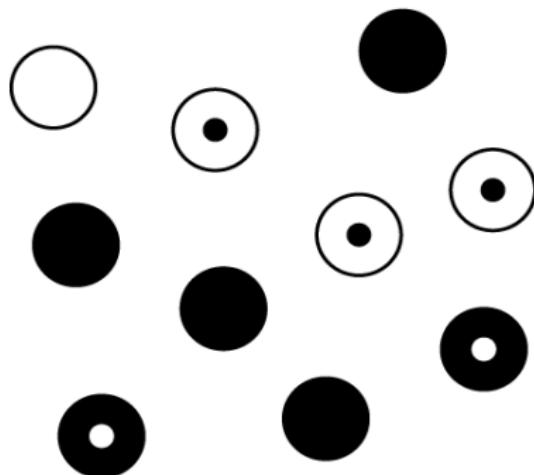
Let's start with joint probability and the simple example that [Paul Lewis](#) gives in his lecture on Bayesian phylogenetics



Slides source: https://molevol.mbl.edu/index.php/Paul_Lewis

Joint probabilities

B = Black S = Solid
W = White D = Dotted



$$\begin{array}{ll} \Pr(B) = 0.6 & \Pr(S) = 0.5 \\ \Pr(W) = 0.4 & \Pr(D) = 0.5 \end{array}$$

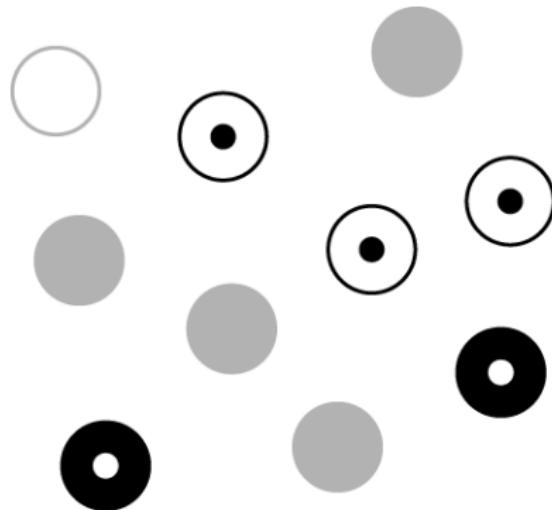
$$\Pr(\bullet) = \Pr(B, D) = 0.2$$

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

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

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

Conditional probabilities



$$\Pr(B|D) = \frac{2}{5} = 0.4$$

Hide all solid marbles
(leaving 5 with dot)

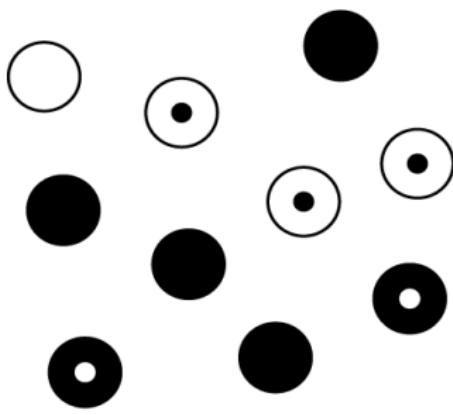
Of those left, 2 are black

Bayes' rule

$\Pr(B, D)$

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

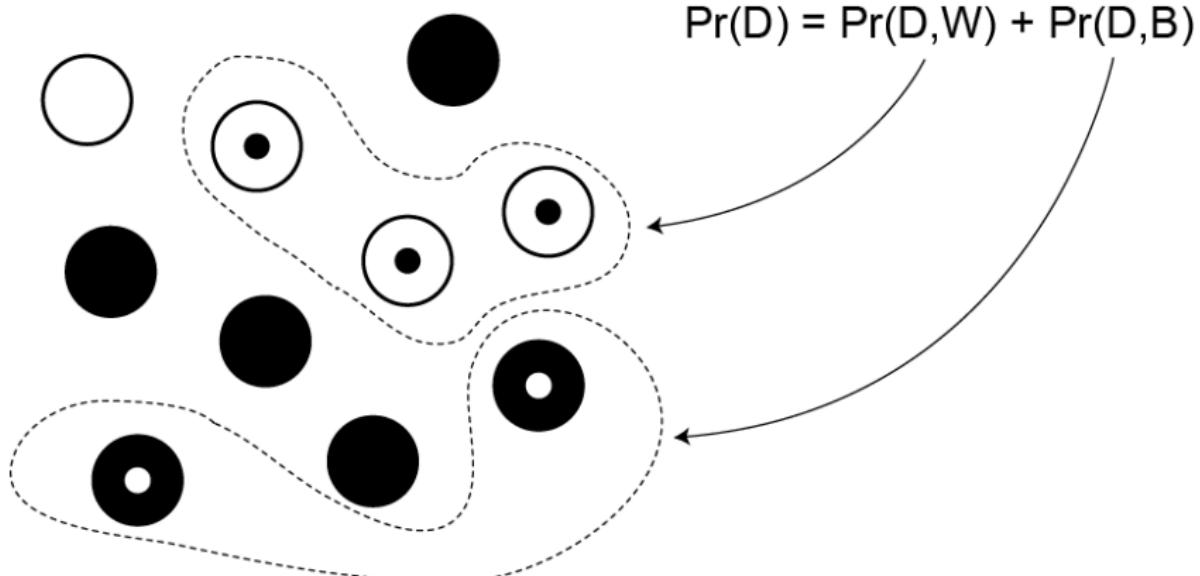
$$\frac{1}{2} \times \frac{2}{5} = \frac{3}{5} \times \frac{1}{3}$$



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

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

Probability of "Dotted"



Bayes' rule (cont.)

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

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

Bayes' rule (cont.)

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

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

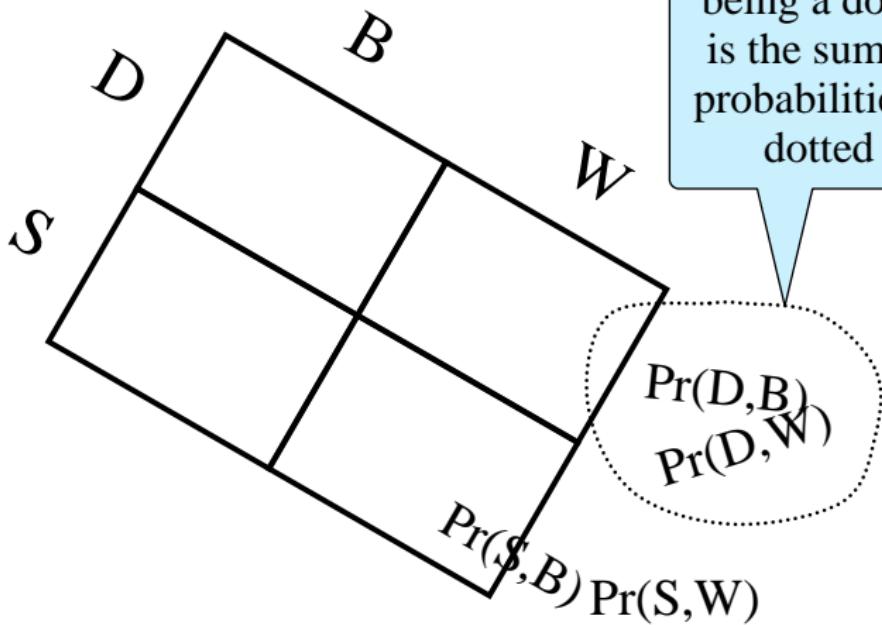
$$\Pr(W|D) = \frac{\Pr(D, W)}{\Pr(D, B) + \Pr(D, W)} \leftarrow \Pr(D)$$

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

Joint probabilities

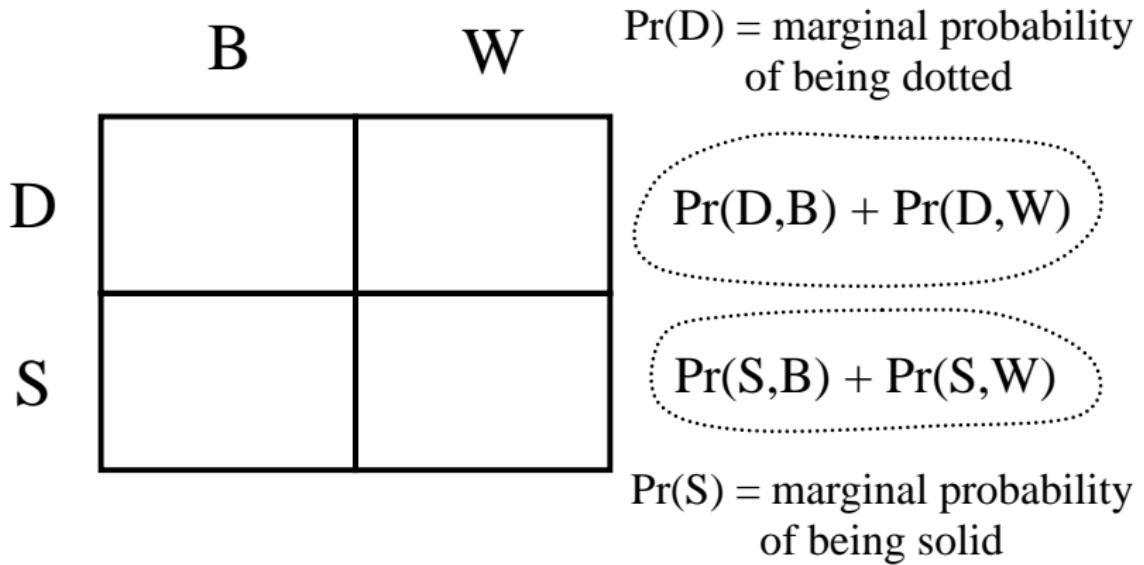
	B	W
D	$\Pr(D,B)$	$\Pr(D,W)$
S	$\Pr(S,B)$	$\Pr(S,W)$

Marginalizing over colors



Marginal probability of being a dotted marble is the sum of all joint probabilities involving dotted marbles

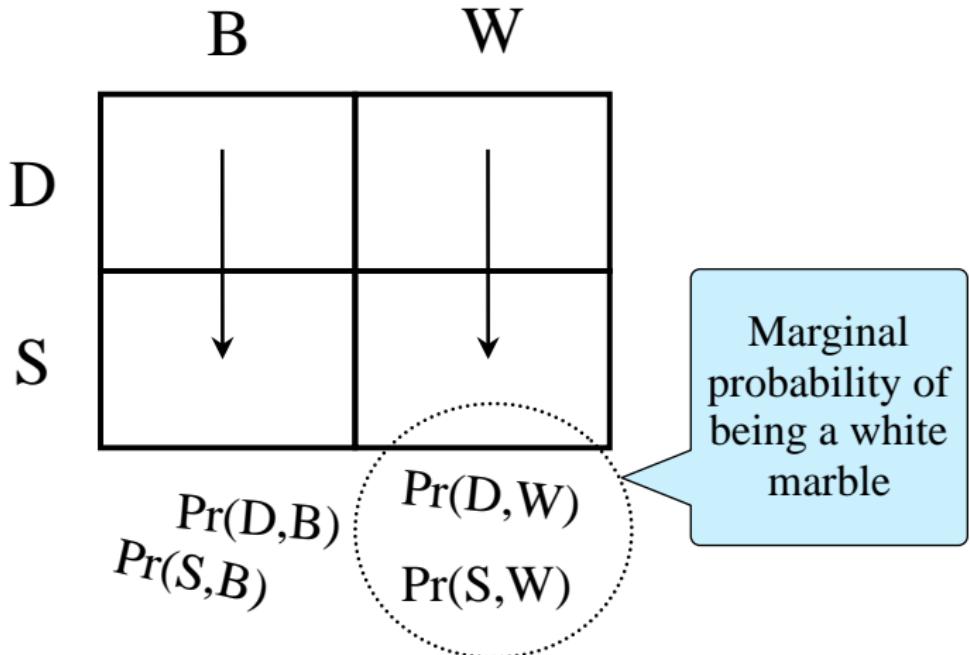
Marginal probabilities



Joint probabilities

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

Marginalizing over "dottedness"



Bayes' rule (cont.)

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

Bayes' rule in Statistics

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

D refers to the "observables" (i.e. the **Data**)

θ refers to one or more "unobservables"
(i.e. **parameters** of a model, or the **model itself**):

- *tree model* (i.e. tree topology)
- *substitution model* (e.g. JC, F84, GTR, etc.)
- *parameter* of a substitution model (e.g. a branch length, a base frequency, transition/transversion rate ratio, etc.)
- *hypothesis* (i.e. a special case of a model)
- a *latent variable* (e.g. ancestral state)

BAYESIAN INFERENCE

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

The probability represents the researcher's **degree of belief**.

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

BAYES' RULE

$$\Pr(\text{Model} \mid \text{Data}) = \frac{\Pr(\text{Data} \mid \text{Model}) \Pr(\text{Model})}{\Pr(\text{Data})}$$

BAYES' RULE

posterior probability

$$\Pr(\text{Model} \mid \text{Data}) = \frac{\Pr(\text{Data} \mid \text{Model}) \Pr(\text{Model})}{\Pr(\text{Data})}$$

BAYES' RULE

$$\Pr(\text{Model} \mid \text{Data}) = \frac{\Pr(\text{Data} \mid \text{Model}) \Pr(\text{Model})}{\Pr(\text{Data})}$$

↓
likelihood

BAYES' RULE

$$\Pr(\text{Model} \mid \text{Data}) = \frac{\Pr(\text{Data} \mid \text{Model}) \Pr(\text{Model})}{\Pr(\text{Data})}$$

↑
prior probability

BAYES' RULE

$$\Pr(\text{Model} \mid \text{Data}) = \frac{\Pr(\text{Data} \mid \text{Model}) \Pr(\text{Model})}{\Pr(\text{Data})}$$

↑
marginal probability of the data

BAYES' RULE

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

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

$\sum_{\delta} \Pr(D | \delta) \Pr(\delta)$ is the likelihood **marginalized** over all possible values of δ .

BAYES' RULE

The posterior probability **density** a continuous parameter θ conditional on the data D is

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

$\int_{\theta} f(D | \theta)f(\theta)d\theta$ is the likelihood **marginalized** over all possible values of θ .

PRIORS

Priors 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

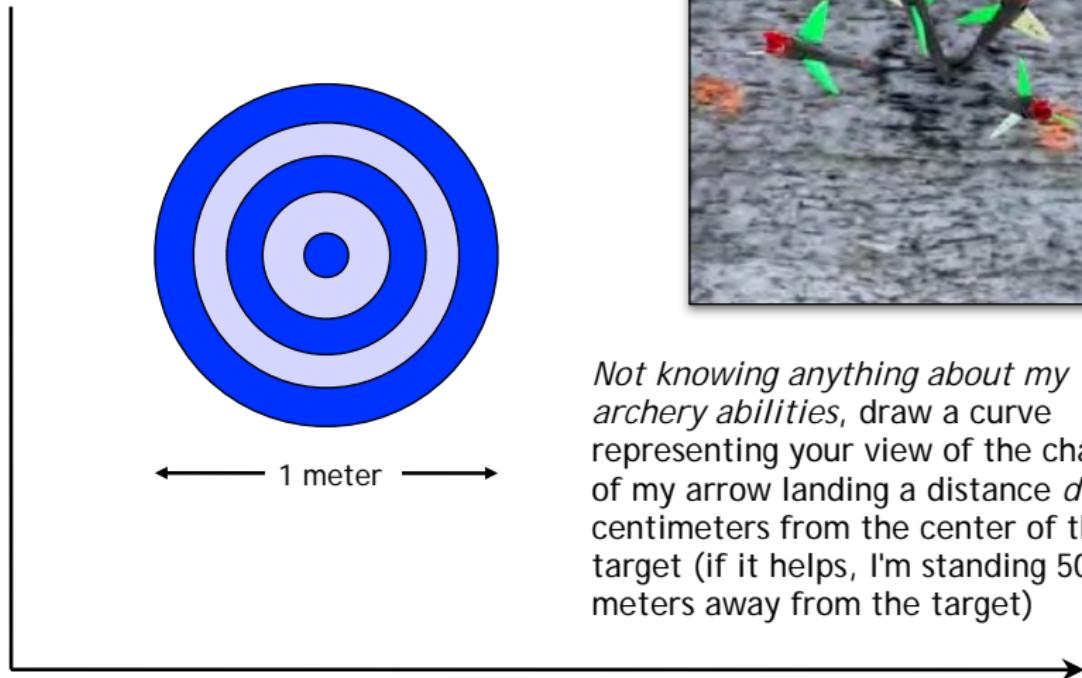
PRIORS

Paul Lewis gives a clear example of a prior in action...



Slides source: https://molevol.mbl.edu/index.php/Paul_Lewis

If you had to guess...



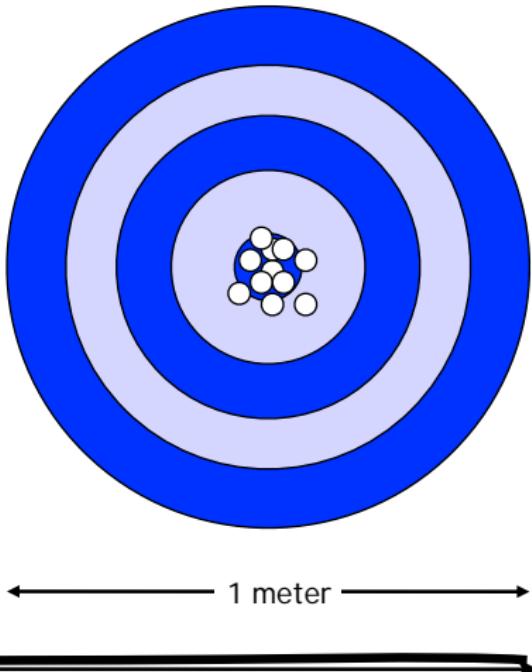
Not knowing anything about my archery abilities, draw a curve representing your view of the chances of my arrow landing a distance d centimeters from the center of the target (if it helps, I'm standing 50 meters away from the target)



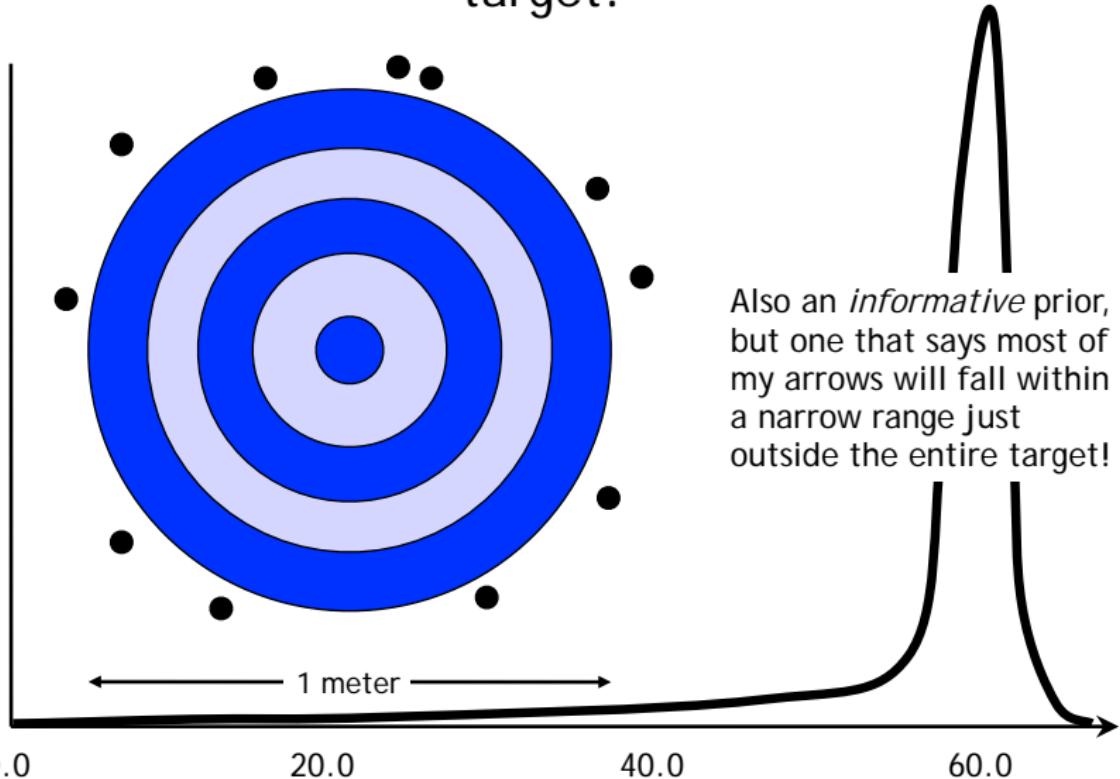
Photo by Tracy Heath

Case 1: assume I have talent

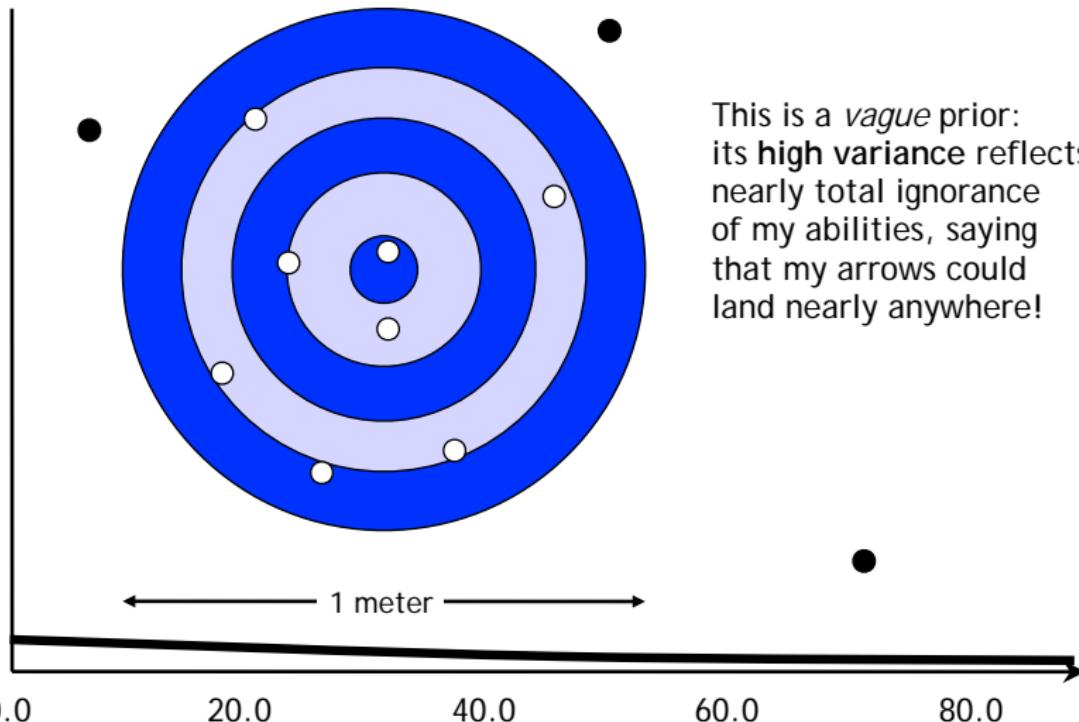
An *informative* prior (low variance) that says most of my arrows will fall within 20 cm of the center (thanks for your confidence!)



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



Case 3: assume I have no talent

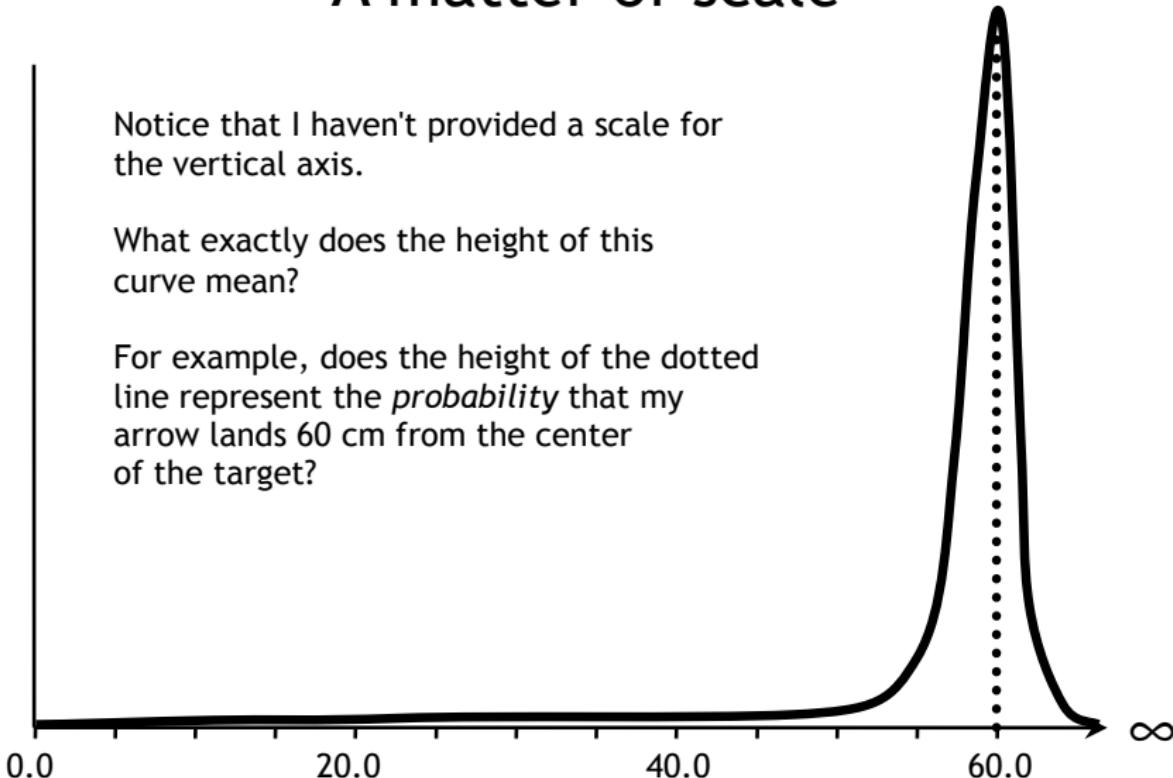


A matter of scale

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?



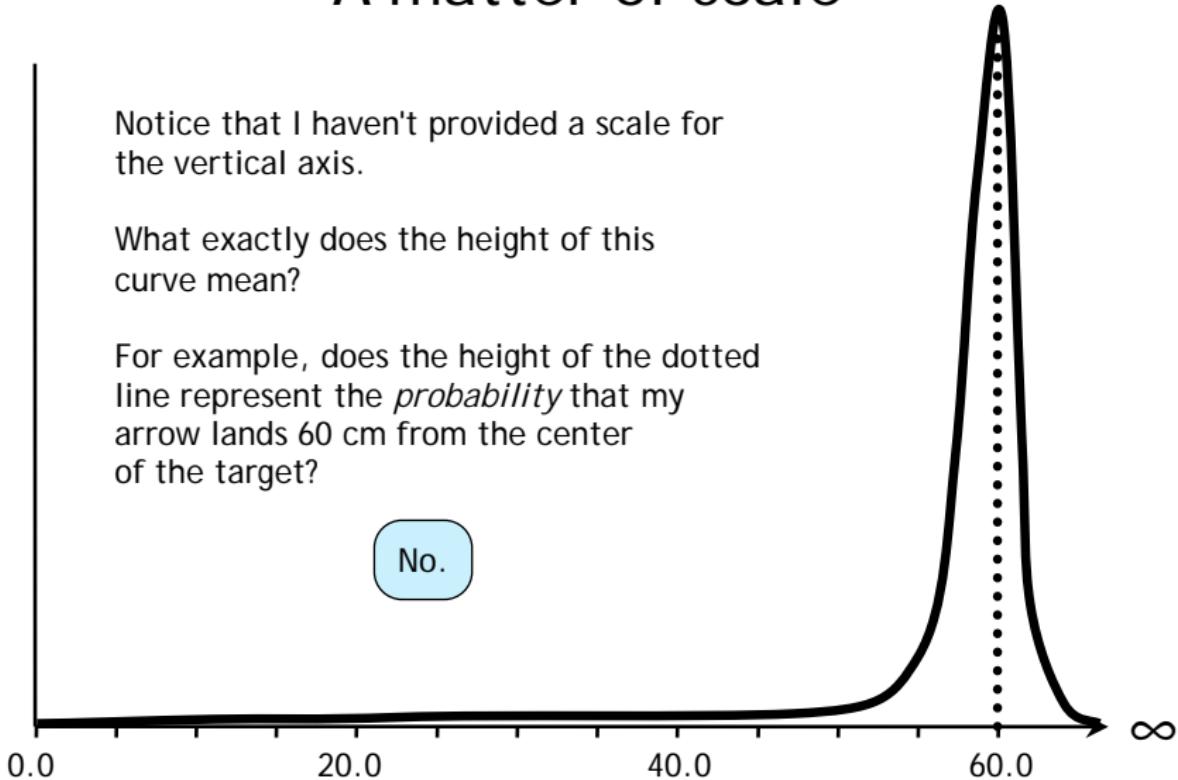
A matter of scale

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?

No.

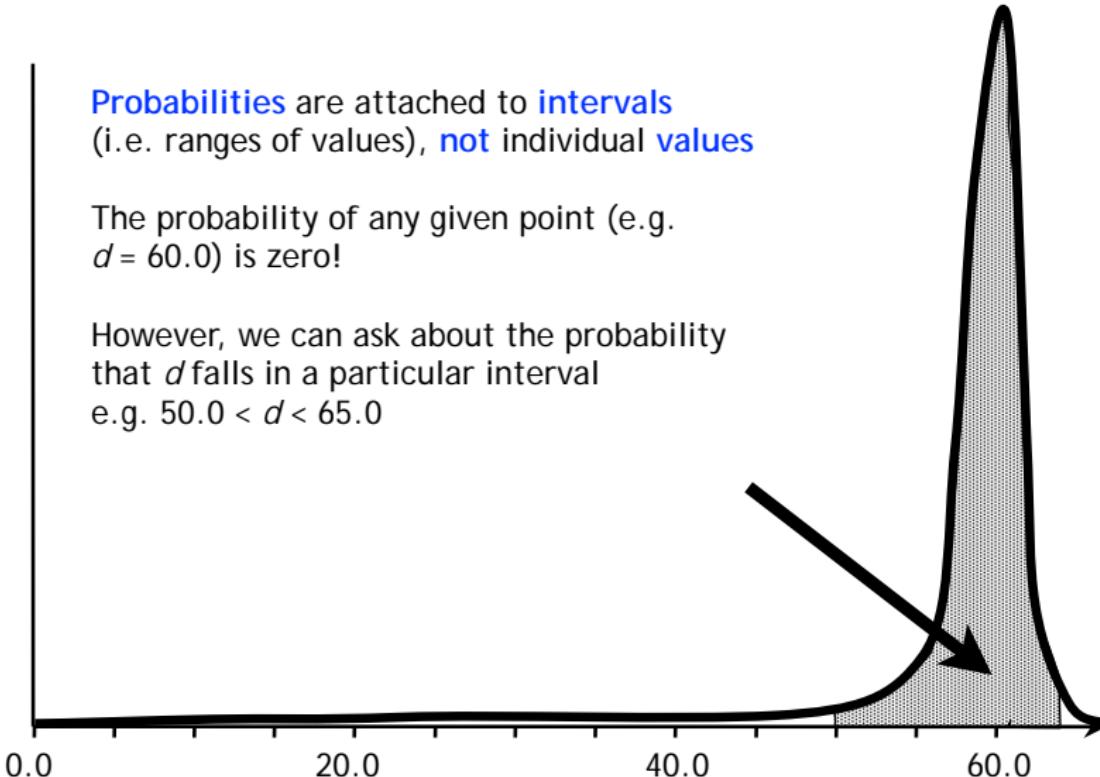


Probabilities are associated with intervals

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. $50.0 < d < 65.0$

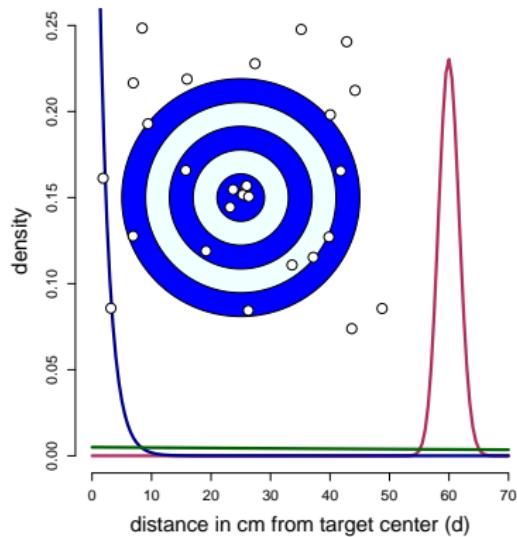


PRIORS: ARCHERY EXAMPLE

Let's continue with the archery example: we may assume a gamma-prior distribution on my archery skill (distance from bullseye = d) with a shape parameter α and a rate parameter β .

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

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



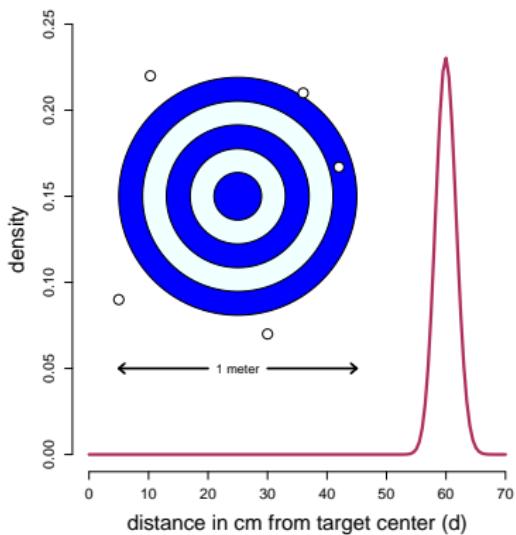
This requires us to assign values for α and β based on our **prior belief**

PRIORS: ARCHERY EXAMPLE

Let's assume that I will consistently miss the target, this corresponds to a gamma distribution with a mean (m) of 60 and a variance (v) of 3.

mean = accuracy

variance = precision



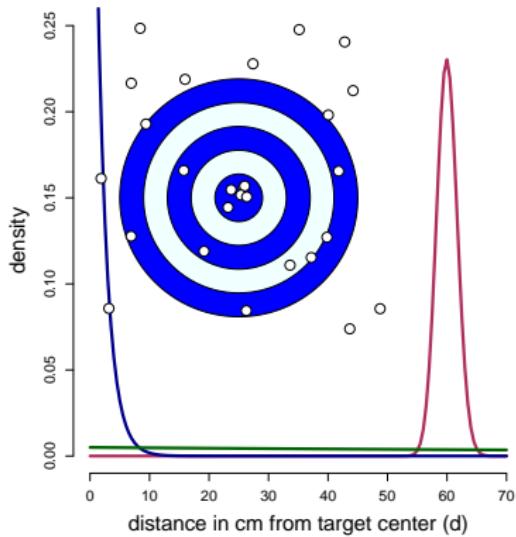
PRIORS: ARCHERY EXAMPLE

If we have some prior knowledge of the mean (m) and variance (v) of the gamma distribution, we can compute α and β .

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

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

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



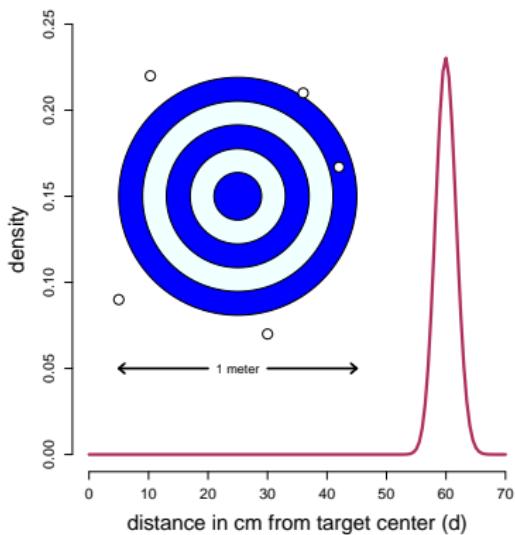
PRIORS: ARCHERY EXAMPLE

Let's assume that I will consistently miss the target, this corresponds to a gamma distribution with a mean (m) of 60 and a variance (v) of 3.

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

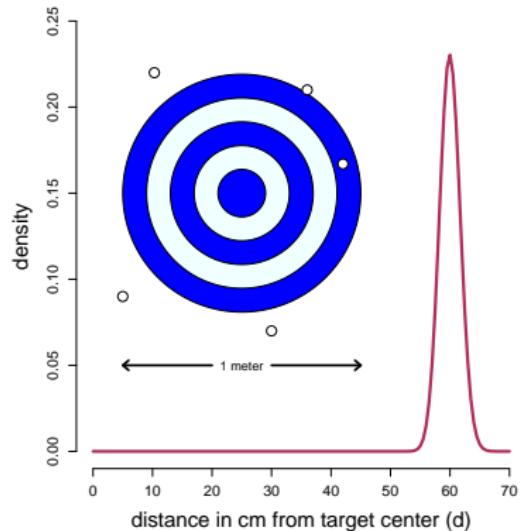
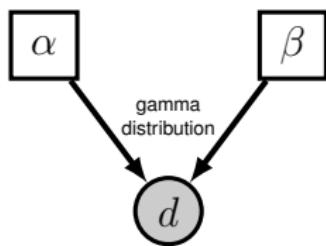
$$\beta = \frac{60}{3} = 20$$

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



PRIORS: ARCHERY EXAMPLE

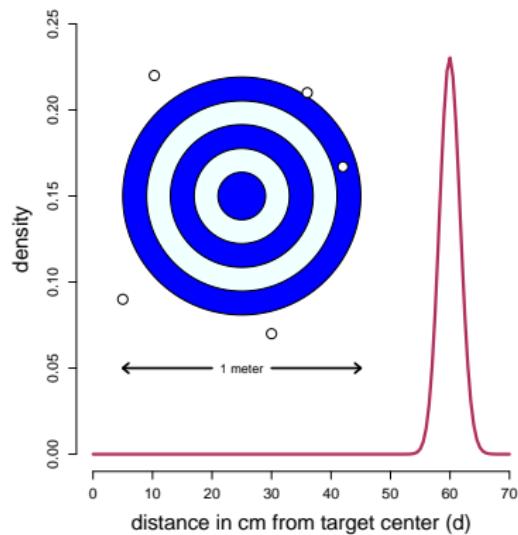
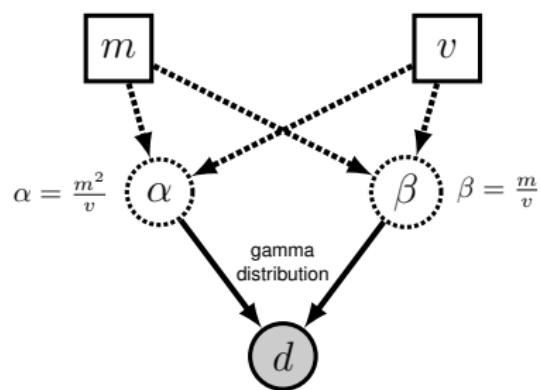
Another way of expressing $d \sim \text{Gamma}(\alpha, \beta)$ is with a **probabilistic graphical model**



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.

PRIORS: ARCHERY EXAMPLE

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

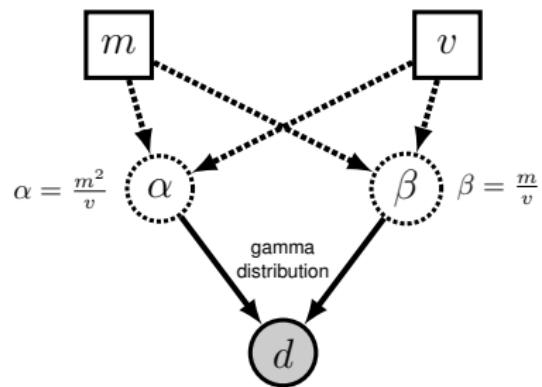


Sometimes it's better to use alternative parameterization.
We may have more intuition about mean and variance than we have about shape and rate.

PRIORS: ARCHERY EXAMPLE

If somehow we happened to **know** the true mean and variance of my accuracy at the archery range, we can easily calculate the likelihood of any observed shot:

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



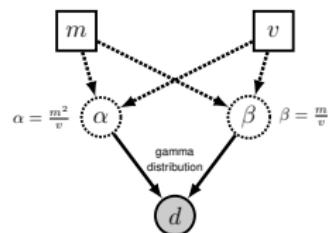
$$f(d = 39.76 | \alpha = 1200, \beta = 20) = 7.89916e - 40$$

RevBayes Demo: ARCHERY Accuracy

RevBayes

Fully integrative Bayesian inference of phylogenetic parameters using probabilistic graphical models and an interpreted language

<http://RevBayes.com>

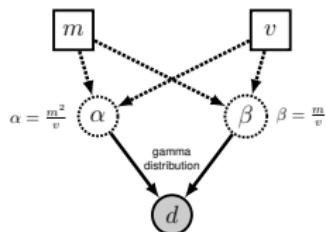


Höhna, Landis, Heath, Boussau, Lartillot, Moore, Huelsenbeck, Ronquist. 2016. RevBayes: Bayesian phylogenetic inference using graphical models and an interactive model-specification language. *Systematic Biology*. (doi: 10.1093/sysbio/syu021)

GRAPHICAL MODELS IN RevBAYES

Graphical models provide tools for visually & computationally representing complex, parameter-rich probabilistic models

We can depict the conditional dependence structure of various parameters and other random variables



Höhna, Heath, Boussau, Landis, Ronquist, Huelsenbeck. 2014.
Probabilistic Graphical Model Representation in Phylogenetics.
Systematic Biology. (doi: 10.1093/sysbio/syu039)

RevBayes Demo: Model on Archery Skill

The Rev language calculating the probability of 1 data observation observed_shot given a mean and variance.

```
mean <- 60
var <- 3

alpha := (mean * mean) / var
beta := mean / var

observed_shot = 39.76

d ~ dnGamma(alpha,beta)
d.clamp(observed_shot)

d.lnProbability()
```

-90.0366

EXAMPLE: MODEL ON ARCHERY SKILL

What if we **do not know** m and v ?

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

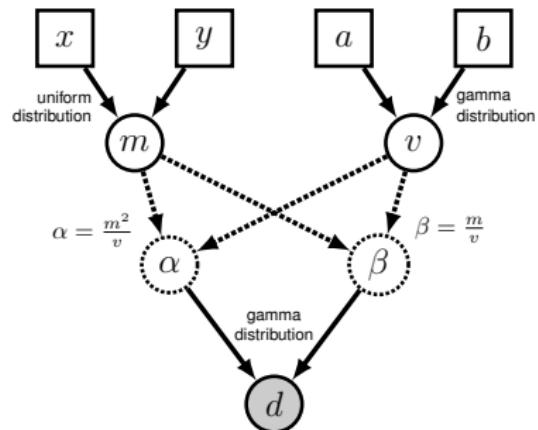
Maximum likelihood methods require us to find the values of m and v that **maximize** $f(d | m, v)$.

Bayesian methods use **prior distributions** to describe our uncertainty in m and v and estimate $f(m, v | d)$.

EXAMPLE: HIERARCHICAL ARCHERY MODEL

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

Now x and y are the parameters of the uniform prior distribution on m and a and b are the shape and rate parameters of the gamma prior distribution on v .



EXAMPLE: HIERARCHICAL ARCHERY MODEL

The values we choose for the parameters of these hyperprior distributions should reflect our prior knowledge. The previous observed shot was **39.76 cm**, thus we may use this to parameterize our hyperpriors for analysis of future observations.

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

$$x = 10.0$$

$$y = 50.0$$

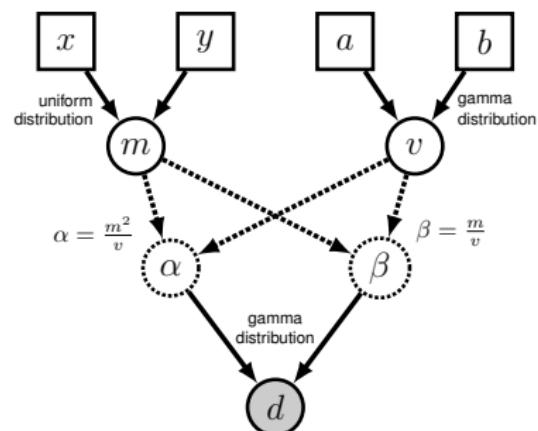
$$\mathbb{E}(m) = 30.0$$

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

$$a = 20.0$$

$$b = 2.0$$

$$\mathbb{E}(v) = 10.0$$



RevBayes Demo: HIERARCHICAL ARCHERY MODEL

The Rev language specifying a hierarchical model on shot accuracy based on 1 new observation.

```
mean ~ dnUnif(10,50)
var ~ dnGamma(20,2)

alpha := (mean * mean) / var
beta := mean / var

observed_shot = 35.21

d ~ dnGamma(alpha,beta)
d.clamp(observed_shot)

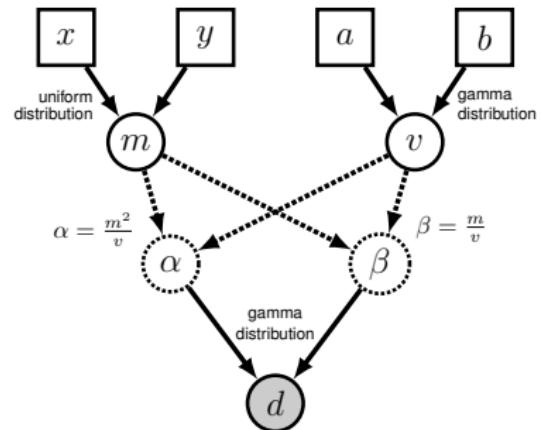
d.lnProbability()
```

depends on initial value of mean & var

EXAMPLE: HIERARCHICAL ARCHERY MODEL

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

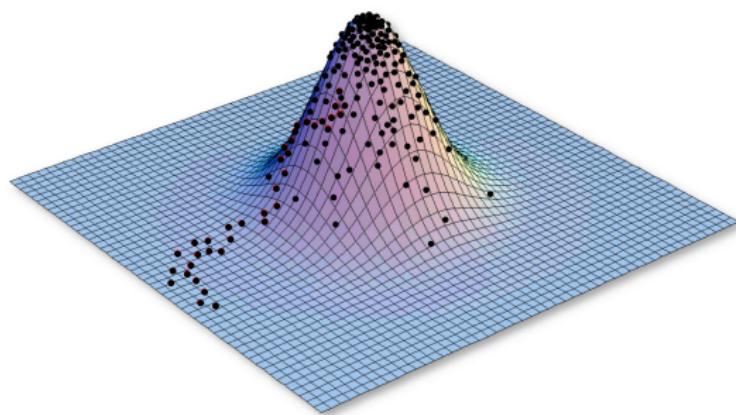
$$\begin{aligned}m &\sim \text{Uniform}(x, y) \\v &\sim \text{Gamma}(a, b) \\d &\sim \text{Gamma}(\alpha, \beta)\end{aligned}$$



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

MARKOV CHAIN MONTE CARLO (MCMC)

An algorithm for approximating the posterior distribution



Metropolis, Rosenbluth, Rosenbluth, Teller, Teller. 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*.

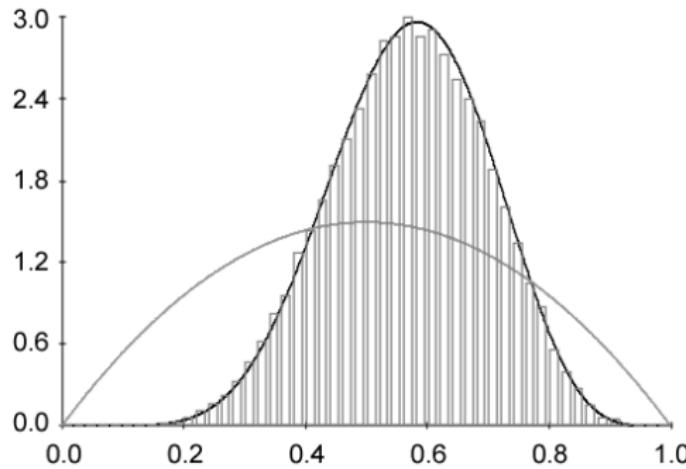
MARKOV CHAIN MONTE CARLO (MCMC)

More on MCMC from Paul Lewis and his lecture on Bayesian phylogenetics



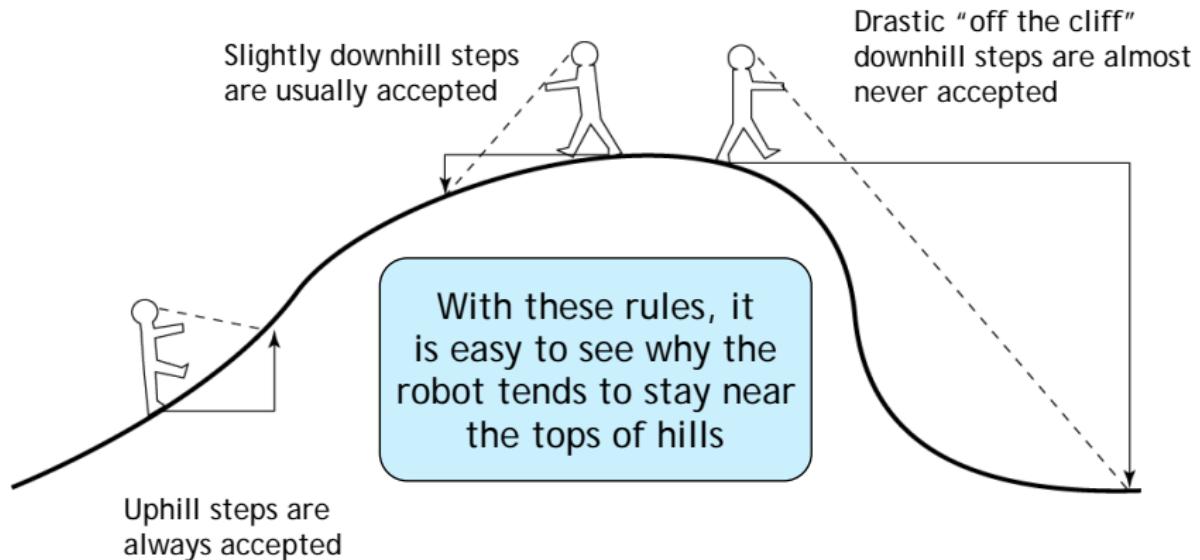
Slides source: https://molevol.mbl.edu/index.php/Paul_Lewis

Markov chain Monte Carlo (MCMC)

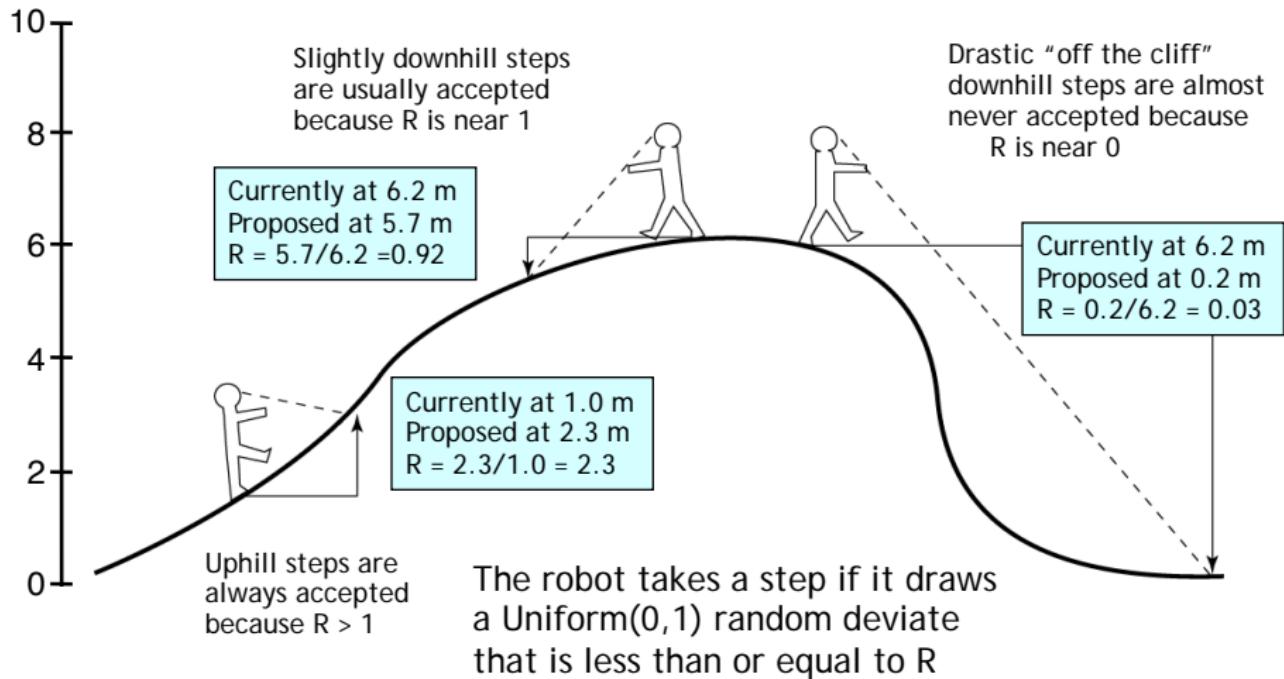


For more complex problems, we might settle for a
good approximation
to the posterior distribution

MCMC robot's rules



(Actual) MCMC robot rules



Cancellation of marginal likelihood

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

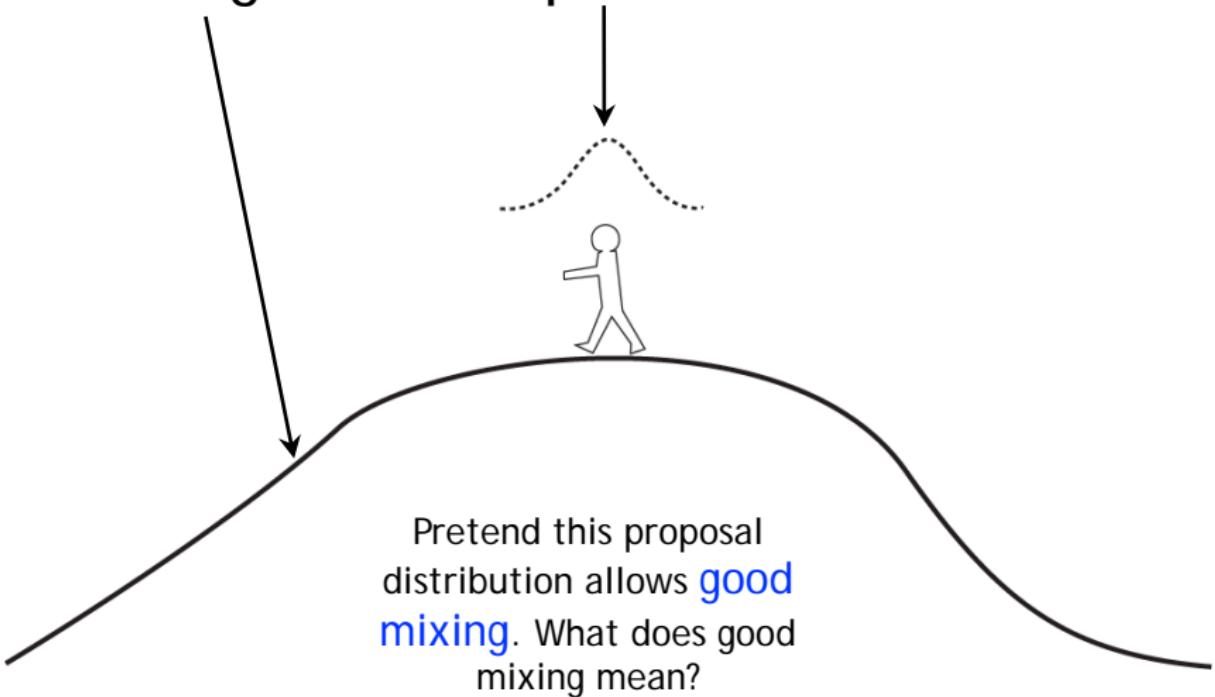
$$\frac{f(\theta^*|D)}{f(\theta|D)} = \frac{\frac{f(D|\theta^*)f(\theta^*)}{\cancel{f(D)}}}{\frac{f(D|\theta)f(\theta)}{\cancel{f(D)}}} = \frac{f(D|\theta^*)f(\theta^*)}{f(D|\theta)f(\theta)}$$

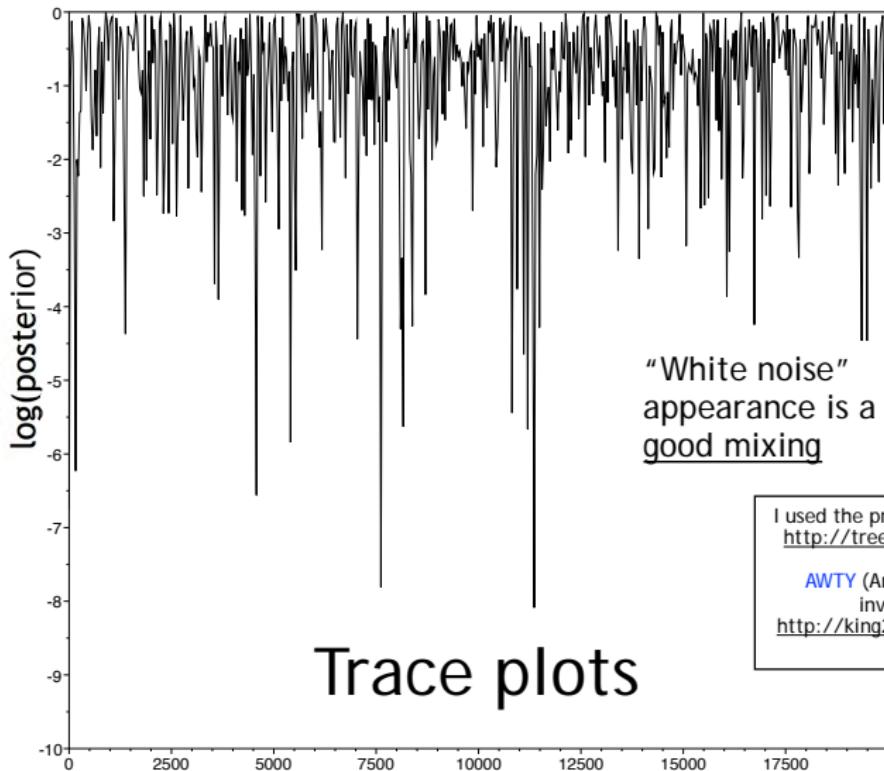
Posterior
odds

Likelihood
ratio

Prior odds

Target vs. Proposal Distributions





Trace plots

I used the program **Tracer** to create this plot:
<http://tree.bio.ed.ac.uk/software/tracer/>

AWTY (Are We There Yet?) is useful for
investigating convergence:
[http://king2.scs.fsu.edu/CEBProjects/awty/
awty_start.php](http://king2.scs.fsu.edu/CEBProjects/awty/awty_start.php)

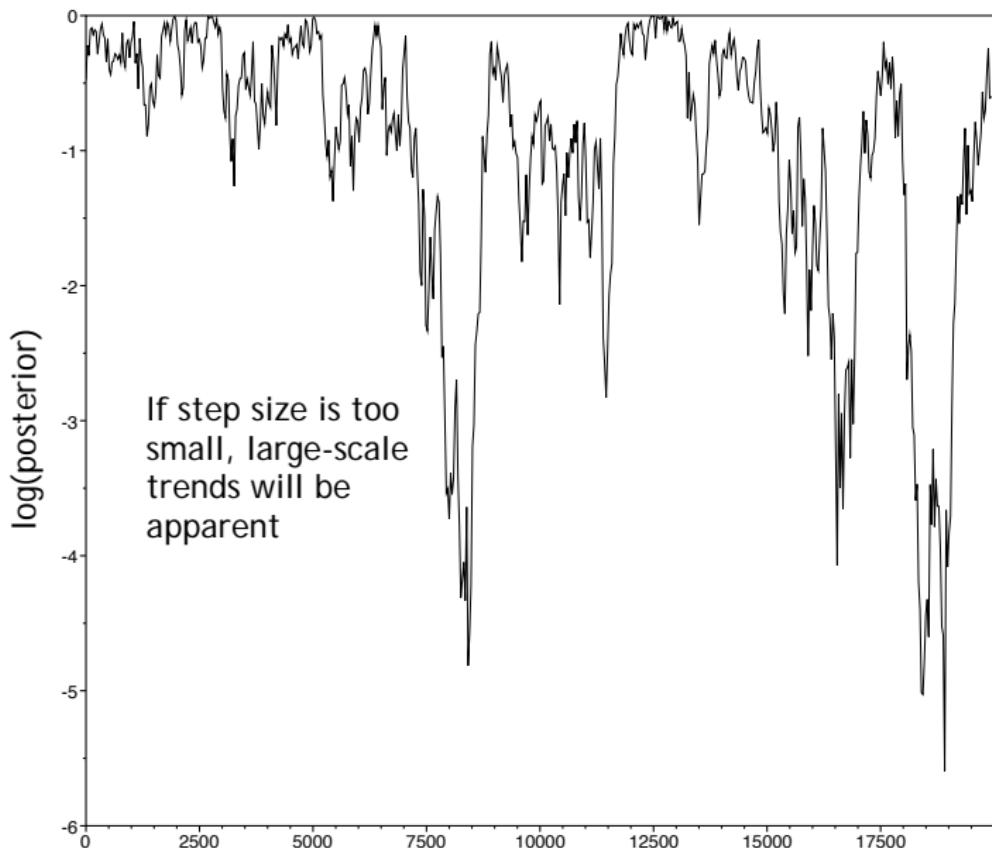
Target vs. Proposal Distributions

Proposal distributions
with **smaller variance**...



Disadvantage: robot takes
smaller steps, more time
required to explore the
same area

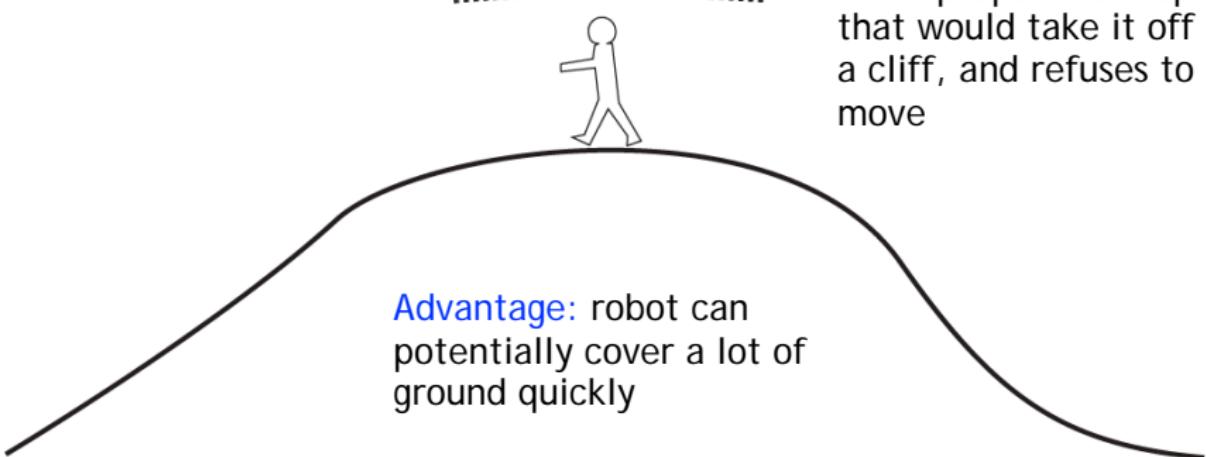
Advantage: robot seldom
refuses to take proposed
steps

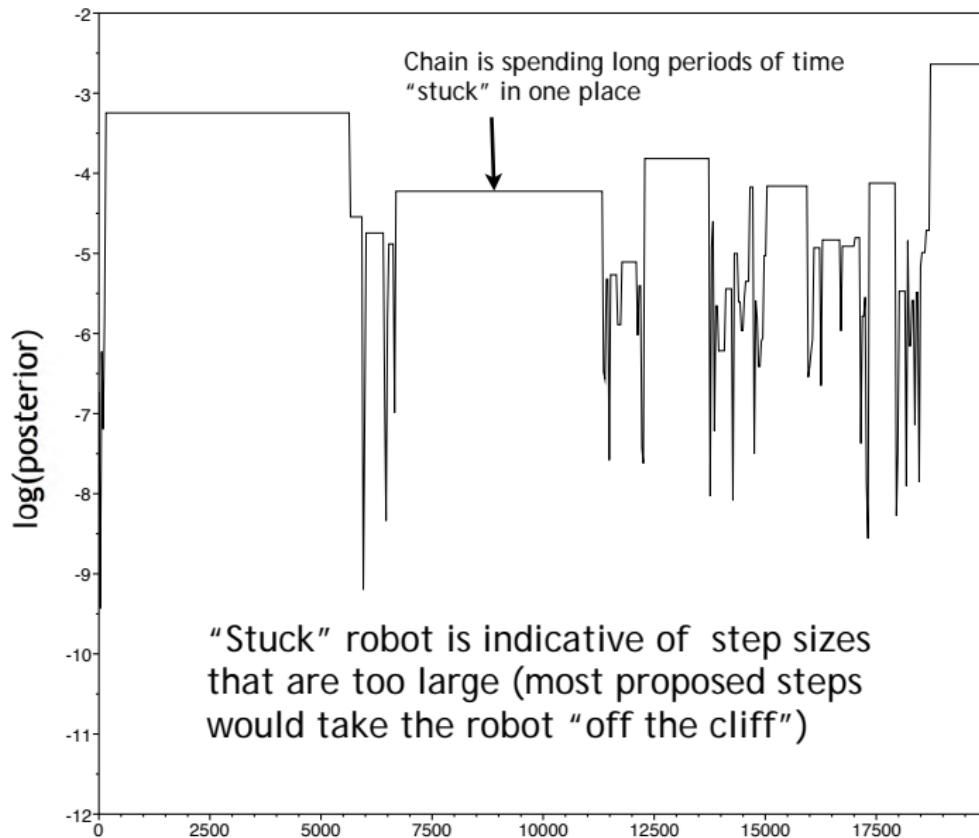


Target vs. Proposal Distributions

Proposal distributions
with **larger variance**...

Disadvantage: robot
often proposes a step
that would take it off
a cliff, and refuses to
move



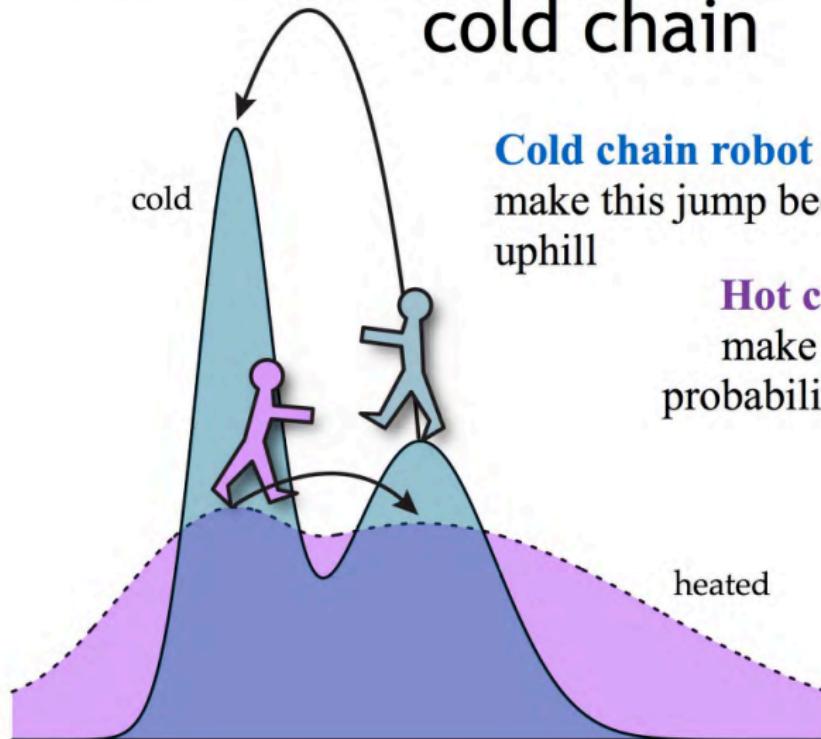


Metropolis-coupled Markov chain Monte Carlo (MCMCMC)

- MCMCMC involves running **several chains simultaneously**
- The **cold chain** is the one that counts, the rest are **heated chains**
- Chain is heated by raising densities to a power less than 1.0 (values closer to 0.0 are warmer)

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

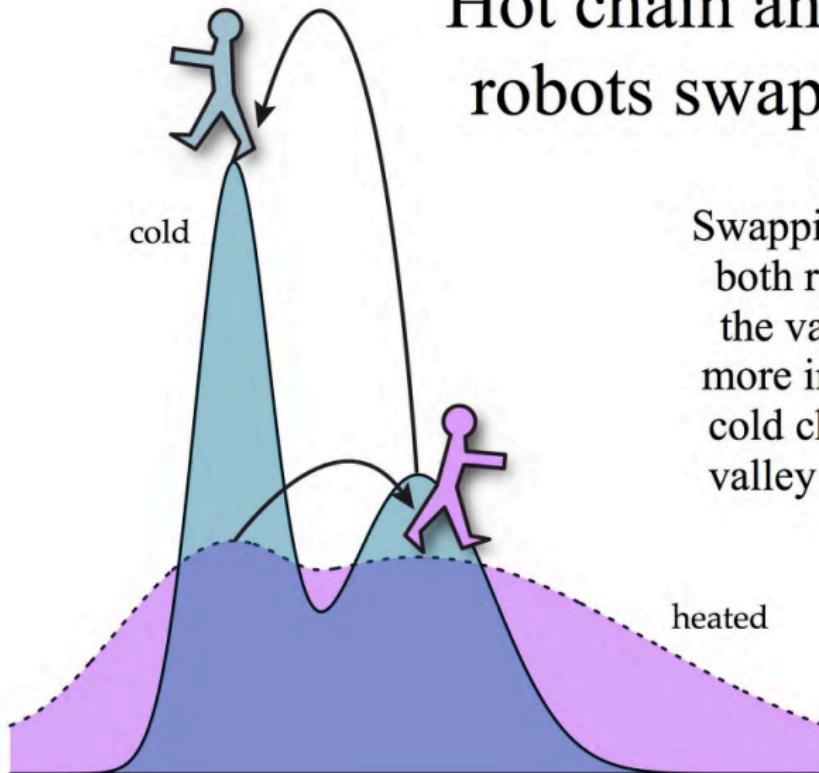
Heated chains act as scouts for the cold chain



Cold 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

Hot chain and cold chain robots swapping places



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

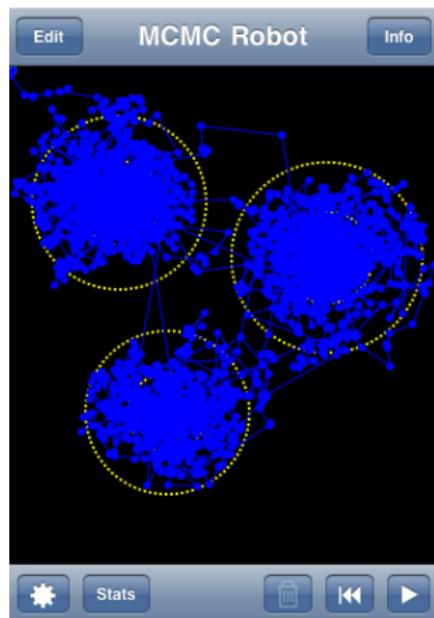
MARKOV CHAIN MONTE CARLO (MCMC)

Thanks, Paul!

Slides source: https://molevol.mbl.edu/index.php/Paul_Lewis

See MCMCRobot, a helpful software program for learning MCMC by Paul Lewis

<http://www.mcmcrobot.org>



RevBayes Demo: HIERARCHICAL ARCHERY MODEL

The Rev language specifying the MCMC sampler for the hierarchical model on archery accuracy.

```
... # model specification from previous demo

mymodel = model(beta)

moves[1] = mvSlide(mean,delta=1.0,tune=true,weight=3.0)
moves[2] = mvScale(var,lambda=1.0,tune=true,weight=3.0)

monitors[1] = mnModel(file="archery_mcmc_1.log",printgen=10, ...)
monitors[2] = mnScreen(printgen=1000, mean, var)

mymcmc = mcmc(mymodel, monitors, moves,nruns=1)

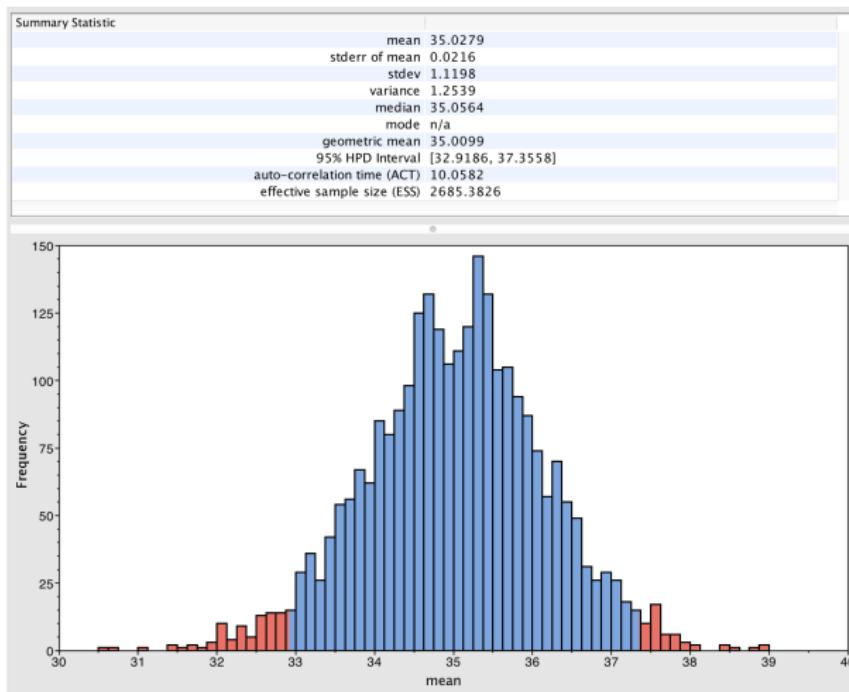
mymcmc.burnin(generations=10000,tuningInterval=1000)

mymcmc.run(generations=40000,underPrior=false)
```

MCMC screen output

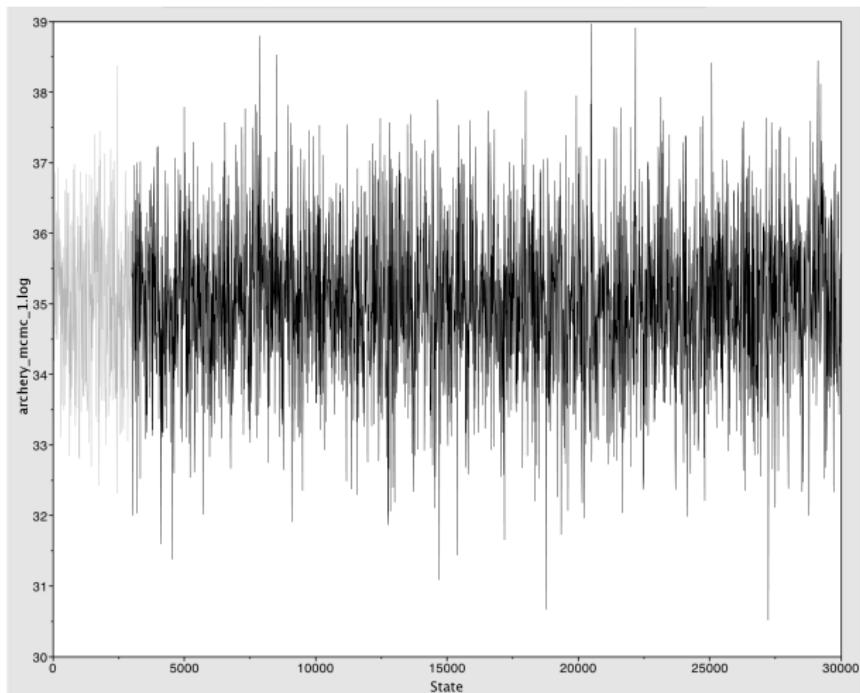
RevBAYES DEMO: HIERARCHICAL ARCHERY MODEL

Summary of the MCMC sample for the mean distance from target center.



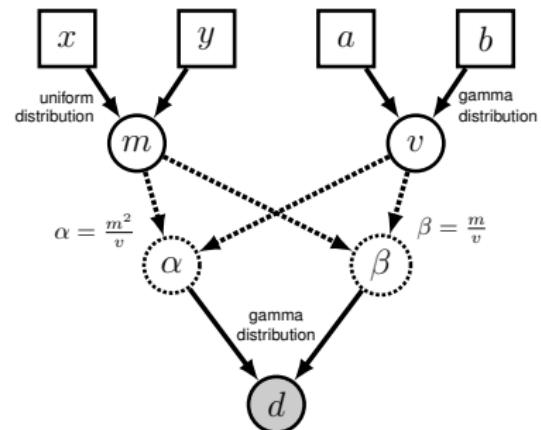
RevBAYES DEMO: HIERARCHICAL ARCHERY MODEL

The trace-plot of the MCMC samples for the mean distance from target center



EXAMPLE: HIERARCHICAL ARCHERY MODEL

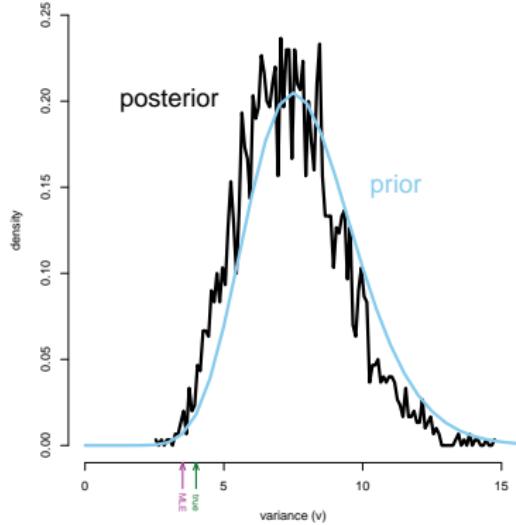
Under this model, we do a good job of estimating the mean, but when judging archery skill, precision (variance) is as (if not more) important than accuracy



Thus, it is also worth evaluating the estimated posterior distribution for the variance component of our model

EXAMPLE: VARIANCE

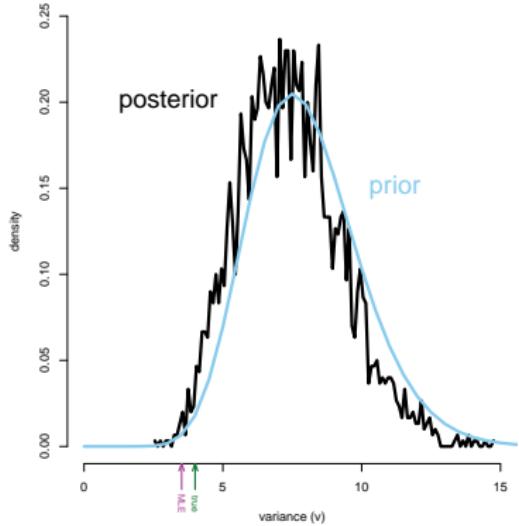
The posterior estimate of the variance (v) is quite different from the true value (4.0) and from the highest likelihood value found by our MCMC (MLE = 3.51374).



This indicates that the prior is having a strong influence on the posterior. Why do you think that is?

EXAMPLE: VARIANCE

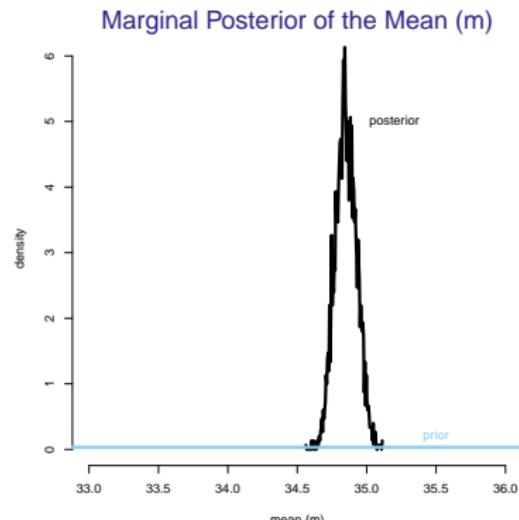
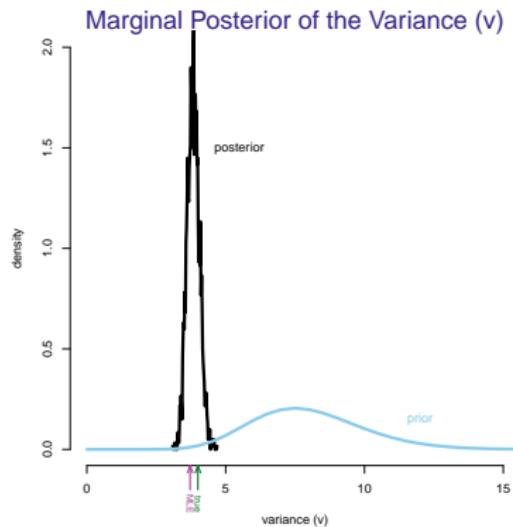
When the prior closely matches the posterior, it can indicate that the data are not very informative for this parameter.



Remember that our data were only 6 observed shots. What would happen if I had 600 arrows?

EXAMPLE: WITH A LOT MORE DATA

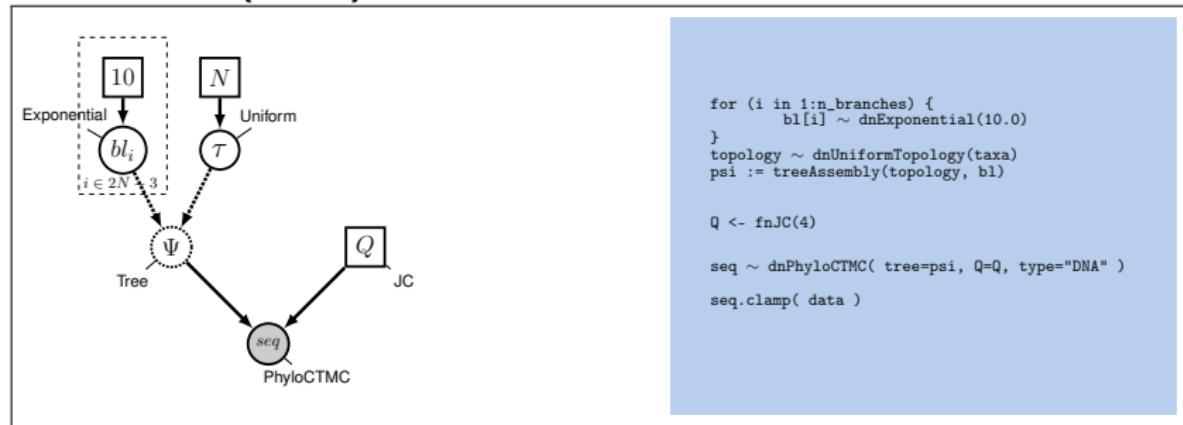
With 100X more observations, we can estimate the mean and variance with greater precision.



BAYESIAN PHYLOGENETICS

How is this applied to phylogenetic inference?

Jukes-Cantor (1969) on an unrooted tree



(image source [RevBayes Substitution Models Tutorial](#))

We can assemble a phylogenetic model in the same way, using previously described models and probability distributions as priors.

BAYESIAN PHYLOGENETICS

With a defined model we simply then have to:

- draw starting values for every random variable in the model
- define moves on each random variable that propose new values
- then for each step in our MCMC, choose a parameter and update it according to the correct proposal.
 - propose a new tree topology and accept or reject
 - propose a new model parameter value and accept or reject
- save the current state of every random variable (tree, branch lengths, base frequencies, etc.) after every k number of states
- after n MCMC steps, evaluate the run for signs of non-convergence
- summarize the tree and other parameters