# **Bayesian Phylogenetics**

Paul O. Lewis Department of Ecology & Evolutionary Biology University of Connecticut

25 January 2013 Workshop on Molecular Evolution Český Krumlov

# An Introduction to Bayesian Phylogenetics

- Bayesian inference in general
- Markov chain Monte Carlo (MCMC)
- Bayesian phylogenetics
- Prior distributions
- Bayesian model selection

#### I. Bayesian inference in general

#### Joint probabilities S $\mathbf{B} =$ Black Solid = White $\mathbf{W} =$ D = Dotted Pr(B) = 0.6 Pr(S) = 0.5Pr(W) = 0.4 Pr(D) = 0.5 $Pr(\mathbf{O}) = Pr(B, D) = 0.2$ $\bullet$ $Pr(\bigcirc) = Pr(B, S) = 0.4$ $Pr(\odot) = Pr(W, D) = 0.3$ $Pr(\bigcirc) = Pr(W, S) = 0.1$

#### Conditional probabilities







## Probability of "Dotted"



## Bayes' rule (cont.)

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

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

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

### Joint probabilities



#### Marginalizing over colors



#### Marginal probabilities





### Bayes' rule (cont.)

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

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**)

 $\theta$  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)

#### Bayes' rule in statistics



#### Simple (albeit silly) paternity example

 $\theta_1$  and  $\theta_2$  are assumed to be the only possible fathers, child has genotype Aa, mother has genotype aa, so child must have received allele **A** from the true father. Note: the data in this case is the child's genotype (Aa)

Possibilities	$\theta_1$	$\theta_2$	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

The prior can be your friend Suppose the test for a **rare** disease is 99% accurate.

$$Pr(+|disease) = 0.99$$
  
 $Pr(+|healthy) = 0.01$   
 $\int$   $\int$  datum hypothesis

Suppose further I test positive for the disease. (Note that we do not need to consider the case of a negative test result.)

It is very tempting to (mis)interpret the likelihood as a posterior probability and conclude "There is a 99% chance that I have the disease."

### The prior can be your friend

The posterior probability is 0.99 only if the **prior probability** of having the disease is 0.5:

$$Pr(disease|+) = \frac{Pr(+|disease)\left(\frac{1}{2}\right)}{Pr(+|disease)\left(\frac{1}{2}\right) + Pr(+|healthy)\left(\frac{1}{2}\right)} \\ = \frac{(0.99)\left(\frac{1}{2}\right)}{(0.99)\left(\frac{1}{2}\right) + (0.01)\left(\frac{1}{2}\right)} = 0.99$$

If, however, the prior odds against having the disease are a million to 1, then the posterior probability is much more reassuring:

$$Pr(disease|+) = \frac{(0.99) \left(\frac{1}{1000000}\right)}{(0.99) \left(\frac{1}{1000000}\right) + (0.01) \left(\frac{999999}{1000000}\right)} \approx 0.0001$$

#### An important caveat

This (rare disease) example involves a **tiny amount of data** (one observation) and an extremely **informative prior**, and gives the impression that maximum likelihood (ML) inference is not very reliable.

However, in phylogenetics, we often have **lots of data** and use much **less informative priors**, so in phylogenetics ML inference is generally **very reliable**.

### Discrete vs. Continuous

- So far, we've been dealing with discrete hypotheses (e.g. either this father or that father, have disease or don't have disease)
- In phylogenetics, substitution models represent an infinite number of hypotheses (each combination of parameter values is in some sense a separate hypothesis)
- How do we use Bayes' rule when our hypotheses form a continuum?

#### Bayes' rule: continuous case



#### 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 from the center of the target (if it helps, I'm standing 50 meters away from the target)

0.0

 $\infty$ 

#### Case 1: assume I have talent



#### 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





#### Densities of various substances

Substance	Density (g/cm <sup>3</sup> )	
Cork	0.24	
Aluminum	2.70	
Gold	19.30	

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

Note: *volume* is appropriate for objects of dimension 3 or higher For 2-dimensions, *area* takes the place of volume For 1-dimension, *linear distance* replaces volume.







# Coin-flipping

y = observed number of heads n = number of flips (sample size) p = (unobserved) proportion of heads

$$\Pr(y|p) = \binom{n}{y} p^y (1-p)^{n-y} = L(p|y)$$

Note that the same formula serves as both the:

- probability of y (if p is fixed)

- likelihood of *p* (if *y* is fixed)

#### Likelihood: why a new term?



#### The posterior is (almost always) more informative than the prior






#### Usually there are many parameters...



marginal probability of the data...

# II. 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



40

#### (Actual) MCMC robot rules



# Cancellation of marginal likelihood

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







#### Target vs. Proposal Distributions





#### Target vs. Proposal Distributions





#### MCRobot (or "MCMC Robot")

Free apps for **Windows** or **iPhone/iPad** available from http://mcmcrobot.org/

Mac version: some day (but see John Huelsenbeck's iMCMC app for MacOS: http://cteg.berkeley.edu/software.html)

Android: some day

#### Tradeoff

- Taking **big steps** helps in jumping from one "island" in the posterior density to another
- Taking **small steps** often results in better mixing
- How can we overcome this tradeoff? **MCMCMC**

# 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 and hot chains swapped



#### Back to MCRobot...

#### The Hastings ratio



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

#### Hastings Ratio



Note that if  $q(\theta|\theta^*) = q(\theta^*|\theta)$ , the Hastings ratio is 1

#### III. Bayesian phylogenetics

# So, what's all this got to do with phylogenetics?



Imagine pulling out trees at random from a barrel. In the barrel, some trees are represented numerous times, while other possible trees are not present. Count 1 each time you see the split separating just A and C from the other taxa, and count 0 otherwise. Dividing by the total trees sampled approximates the true proportion of that split in the barrel.

#### Moving through treespace



## Moving through parameter space



Using  $\kappa$  (ratio of the transition rate to the transversion rate) as an example of a model parameter.

Proposal distribution is the uniform distribution on the interval ( $\kappa$ -d,  $\kappa$ +d)

The "step size" of the MCMC robot is defined by d: a larger d means that the robot will attempt to make larger jumps on average.

### Putting it all together

- Start with random tree and arbitrary initial values for branch lengths and model parameters
- Each generation consists of one of these (chosen at random):
  - Propose a new tree (e.g. Larget-Simon move) and either accept or reject the move
  - Propose (and either accept or reject) a new model parameter value
- Every k generations, save tree topology, branch lengths and all model parameters (i.e. sample the chain)
- After *n* generations, **summarize sample** using histograms, means, credible intervals, etc.

## Marginal Posterior Distribution of $\kappa$



Histogram created from a sample of 1000 kappa values.

Data from Lewis, L., and Flechtner, V. 2002. Taxon 51: 443-451.

#### **IV. Prior distributions**

#### **Common Priors**

- **Discrete uniform** for topologies – exceptions becoming more common
- Beta for proportions
- Gamma or Log-normal for branch lengths and other parameters with support  $[0,\infty)$ 
  - Exponential is common special case of the gamma distribution
- **Dirichlet** for state frequencies and GTR relative rates

#### **Discrete Uniform** distribution for **topologies**

























# Yule model provides joint prior for both topology and divergence times



The rate of speciation under the Yule model ( $\lambda$ ) is constant and applies equally and independently to each lineage. Thus, speciation events get closer together in time as the tree grows because more lineages are available to speciate.

#### Gamma(*a*,*b*) distributions



\*Note: be aware that in many papers the Gamma distribution is defined such that the second (scale) parameter is the *inverse* of the value *b* used in this slide! In this case, the mean and variance would be a/b and  $a/b^2$ , respectively.

Log-normal distribution



**Important:**  $\mu$  and  $\sigma$  do **not** represent the mean and standard deviation of X: they are the mean and standard deviation of  $\log(X)$ !

To choose  $\mu$  and  $\sigma$  to yield a particular mean (*m*) and variance (*v*) for X, use these formulas:  $\log(v + m^2) - \log(m^2)$ 

$$\mu = \log(m^2) - \log(m) - \frac{\log(v + m^2) - \log(v)}{2}$$
  
$$\sigma^2 = \log(v + m^2) - \log(m^2)$$

#### Beta(a,b) gallery



#### Dirichlet(a,b,c,d) distribution

Used for nucleotide relative frequencies:



(stereo pairs)

Flat prior:

a = b = c = d = 1

(no scenario discouraged)

Informative prior: a = b = c = d = 300

(equal frequencies strongly encouraged)

Dirichlet(a,b,c,d,e,f) used for GTR exchangeability parameters.

(Thanks to Mark Holder for suggesting the use of a tetrahedron)

#### **Prior Miscellany**

- priors as rubber bands <
- running on empty
- hierarchical models
- empirical bayes




#### Example: Internal Branch Length Priors



Separate priors applied to internal and external branches

External branch length prior is exponential with mean 0.1

Internal branch length prior is exponential with mean 0.1

This is a reasonably vague internal branch length prior



#### Internal branch length prior mean 0.01

(external branch length prior mean always 0.1)



### Internal branch length prior mean 0.001



## Internal branch length prior mean 0.0001

40 Cyanophora paradoxa – 39 Nephroselmis olivacea - 38 Pteromonas angulos 36 Chlamydomonas reinhardtii - 37 Paulschulzia pseudovolvox — 35 Volvox carteri - 33 Mesostigma viride └─ 34 Mesostigma viride NIES - 32 Chlorokybus atmosphyticus - 20 Chaet oval — 19 Chaet globosum SAG2698 – 25 Mesotaenium caldariorum – 27 Mougeotia sp 758 - 23 Gonatozygon monotaenium - 21 Onychonema sp - 22 Cosmocladium perissum — 31 Entransia fimbriata — 24 Spirogyra maxima 2495 - 26 Zygnema peliosporum 45 - 28 Klebsormidium flaccidum - 29 Klebsormidium subtilissimum - 30 Klebsormidium nitens 16 Coleochaete soluta 32d1 - 15 Coleochaete orbicularis - 17 Coleochaete irregularis – 18 Coleochaete sieminskiana — 14 Tolypella int prolifera — 13 Nitella opaca - 9 Chara connivens - 10 Lamprothamnium macropogon - 11 Lychnothamnus barbatus └ 12 Nitellopsis obtusa - 8 Marchantia polymorpha - 7 Anthoceros formosae - 6 Sphagnum palustre - 3 Huperzia lucidula – 4 Psilotum nudum — 5 Dicksonia antarctica – 2 Taxus baccata 1 Arabidopsis thaliana

## Internal branch length prior mean 0.00001

40 Cyanophora paradoxa – 39 Nephroselmis olivacea - 38 Pteromonas angulos - 37 Paulschulzia pseudovolvox - 35 Volvox carteri – 36 Chlamvdomonas reinhardtii - 34 Mesostigma viride NIES - 33 Mesostigma viride - 32 Chlorokybus atmosphyticus - 23 Gonatozygon monotaenium 22 Cosmocladium perissum - 30 Klebsormidium nitens – 29 Klebsormidium subtilissimum - 21 Onychonema sp — 27 Mougeotia sp 758 — 24 Spirogyra maxima 2495 - 26 Zygnema peliosporum 45 – 25 Mesotaenium caldariorum - 31 Entransia fimbriata - 19 Chaet globosum SAG2698 - 20 Chaet oval - 28 Klebsormidium flaccidum — 15 Coleochaete orbicularis - 17 Coleochaete irregularis - 16 Coleochaete soluta 32d1 – 18 Coleochaete sieminskiana — 13 Nitella opaca - 14 Tolypella int prolifera - 12 Nitellopsis obtusa — 11 Lychnothamnus barbatus - 9 Chara connivens – 10 Lamprothamnium macropogon - 8 Marchantia polymorpha 7 Anthoceros formosae - 3 Huperzia lucidula - 6 Sphagnum palustre - 4 Psilotum nudum – 5 Dicksonia antarctica - 2 Taxus baccata 1 Arabidopsis thaliana

0.1

Internal branch length prior mean 0.000001

The internal branch length prior is calling the shots now, and the likelihood must obey.

#### **Prior Miscellany**

- priors as rubber bands
- running on empty



- hierarchical models
- empirical bayes

#### Running on empty



#### **Prior Miscellany**

- priors as rubber bands
- running on empty
- hierarchical models
- empirical bayes



In a **non-hierarchical** model, all parameters are present in the likelihood function

Prior: Exponential, mean=0.1



# **Hierarchical** models add *hyperparameters* not present in the likelihood function

 $\mu$  is a *hyperparameter* governing the mean of the edge length prior



Prior: Exponential, mean  $\mu$ 

$$L_{k} = \frac{1}{4} \left[ \frac{1}{4} + \frac{3}{4} e^{-4v_{1}/3} \right] \left[ \frac{1}{4} + \frac{3}{4} e^{-4v_{2}/3} \right] \left[ \frac{1}{4} - \frac{1}{4} e^{-4v_{3}/3} \right] \left[ \frac{1}{4} - \frac{1}{4} e^{-4v_{4}/3} \right] \left[ \frac{1}{4} + \frac{3}{4} e^{-4v_{5}/3} \right]$$

During an MCMC analysis,  $\mu$  will hover around a reasonable value, sparing you from having to decide what value is appropriate. You still have to specify a hyperprior, however.

For example, see Suchard, Weiss and Sinsheimer. 2001. MBE 18(6): 1001-1013.

#### **Prior Miscellany**

- priors as rubber bands
- running on empty
- hierarchical models
- empirical bayes



#### **Empirical Bayes**

Empirical Bayes uses the data to determine some aspects of the prior, such as the prior mean. This uses the data twice, which is not acceptable to Bayesian purists

An empirical Bayesian would use the maximum likelihood estimate (MLE) of the length of an average branch here

Prior: Exponential, mean=MLE

$$L_{k} = \frac{1}{4} \left[ \frac{1}{4} + \frac{3}{4} e^{-4v_{1}/3} \right] \left[ \frac{1}{4} + \frac{3}{4} e^{-4v_{2}/3} \right] \left[ \frac{1}{4} - \frac{1}{4} e^{-4v_{3}/3} \right] \left[ \frac{1}{4} - \frac{1}{4} e^{-4v_{4}/3} \right] \left[ \frac{1}{4} + \frac{3}{4} e^{-4v_{5}/3} \right]$$

#### V. Bayesian model selection

AIC is not Bayesian. Why?

 $AIC = 2k - 2\log(\max_{\uparrow} L)$ 

number of free (estimated) parameters maximized log likelihood

AIC is not Bayesian because the **prior is not considered** (and the prior is an important component of a Bayesian model)

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

The **marginal likelihood** (denominator in Bayes' Rule) is commonly used for Bayesian model selection

Represents the (weighted) **average fit of the model** to the observed data (weights provided by the prior)

#### An evolutionary distance example



- Let's compare models JC69 vs. K80
- Parameters:
  - v is edge length (expected no. substitutions/site)
    free in both JC69 and K80 models
  - κ is transition/transversion rate ratio
    - free in K80, set to 1.0 in JC69

#### Likelihood Surface when K80 true



### Likelihood Surface when JC true

Based on simulated data:



#### Harmonic mean method

$$\widehat{f(D|M)} = \frac{n}{\frac{1}{L^{(1)}} + \frac{1}{L^{(2)}} + \dots + \frac{1}{L^{(n)}}} \begin{bmatrix} L^{(i)} = \text{Likelihood (not log-likelihood) calculated for the ith sample from the MCMC analysis} \\ \log BF_{12} = \log \left(\frac{f(D|M_1)}{f(D|M_2)}\right) \\ = \log f(D|M_1) - \log f(D|M_2) \\ \widehat{\uparrow} \end{bmatrix}$$

Most Bayesian programs provide the log of the harmonic mean of the sampled likelihoods for each model you run, so all you need to do is subtract.

Example: MrBayes output	Run	Arithmetic mean	Harmonic mean
	1 2	-22913.52 -22913.52	-22923.02 -22922.68
	TOTAL	-22913.52	-22922.86

*Warning:* the harmonic mean method is **strongly biased** and **should not be used** if more accurate methods are available

Newton, M. A. and A. E. Raftery. 1994. Approximate Bayesian inference with the weighted likelihood bootstrap (with discussion). J. Roy. Stat. Soc. B 56:3–48.

#### Another approach





### How many "stepping stones" (i.e. ratios) are needed?



#### Is steppingstone sampling accurate?



#### How about the harmonic mean method?



