Phylogenetics Paul Lewis lecture series Q & A RL-V3 MPP

Rachel Warnock 08.07.25



Projects

Written report 5 pages max, inc. references and figures Focus mainly on the methods and results

Brief intro & discussion (paragraph each)

Provide access to the project code

New final deadline 31.08.25



Trouble shooting

If you have any questions, please send them via slack over the next few weeks (not email)

Note I will not be available to answer questions July 12 - July 28



Q&A <u>Phylogenetics primer part 1</u> by Paul Lewis

(the answers provided here are my interpretation of these concepts – answers may vary!)



1. In your own words, how would you describe the terms conditional probability and likelihood?

Conditional probability \rightarrow the probability of an event, dependent on the value of some other event

Likelihood \rightarrow the probability of our observations given a set of assumptions (i.e., the model) and parameter values



2. Why do we calculate likelihoods on a log scale?

→ Because likelihoods can get incredibly small — see for yourself using R!



3. How does the probability of transitioning from one character state to another (e.g., from A to T) change over time?

use the Transition Probability tool to explore this further

→ the probability of change increases with time



Felsenstein's pruning algorithm

4a. What do we need this for?

 \rightarrow to calculate the likelihood of a tree (given an alignment and a substitution) model, taking into account all possible ancestral states at every node)

4b. Can you describe the gist of Felsenstein's pruning algorithm?

For a good description of Felsenstein's pruning algorithm see <u>Section 8.8</u> of Phylogenetic Comparative Methods by Harmon



Q&A <u>Phylogenetics primer part 2</u> by Paul Lewis

(the answers provided here are my interpretation of these concepts – answers may vary!)



1. What are the assumptions of the following substitution models? Consider the rate of change between **character states** and the **state frequencies**.

 $JC \rightarrow$ $HKY \rightarrow$

- the rate of change between character states and the state frequencies.
 - $JC \rightarrow$ equal frequencies, equal rates
- $GTR \rightarrow$ unequal frequencies, unequal rates

transversions: $A \leftrightarrow T$ or $G \leftrightarrow C$, transitions: $A \leftrightarrow G$ or $C \leftrightarrow T$

1. What are the assumptions of the following substitution models? Consider

HKY → unequal frequencies, unequal rates between <u>transversions</u> & <u>transitions</u>



variation among characters?

2a. Site specific rates

set of parameters

2b. Invariant sites model

 \rightarrow assign a subset of sites to a "constant" (i.e., non-variable) category

2. Can you briefly describe the following approaches to account for rate

\rightarrow assign sites to separate partitions and allow each partition to have its own



2. Can you briefly describe the follow variation among characters?

2c. Discrete Gamma model

→ calculate the likelihood assuming there are discrete rate categories, e.g., 4. Variation in rate categories is represented by a gamma distribution and the parameters of the gamma distribution are calculated from the data

2. Can you briefly describe the following approaches to account for rate



3. For a given substitution model, what do the values in the Q matrix and the **P** matrix represent?

of change between each character state combination

the branch lengths

- \rightarrow the Q matrix is the instantaneous rate matrix, i.e., the instantaneous rates
- \rightarrow the P matrix is the transition probability matrix, i.e., the probabilities of change between character states after relative time t, which is represented by



Q&A Phylogenetics primer part 3a by Paul Lewis

(the answers provided here are my interpretation of these concepts – answers may vary!)



1. In your own words can you describe each component of Bayes' rule? Which parts are difficult to understand?



Recap



Pr(data | model) Pr(model)

Pr(model | data) =



Likelihood

Pr(model | data) =

The probability of the data given the model assumptions and parameter values

Pr(data | model) Pr(model)



Pr(data | model) Pr(model)

Pr(model | data) =

This represents our prior knowledge of the model parameters

Priors



Pr(data | model) Pr(model)

Pr(model | data) =



Pr(data)

Marginal probability

The probability of the data, given all possible parameter values. Can be thought of as a normalising constant



Reflects our combined knowledge based on the likelihood and the priors

posterior

Pr(data | model) Pr(model)

Pr(model | data) =



Components used to infer trees without considering time

0101... 1101... 0100...

data sequences or characters topology and branch lengths





substitution model



Bayesian tree inference

posterior





1. In your own words can you describe each component of Bayes' rule? Which parts are difficult to understand?

Posterior ~ Likelihood × Priors

The posterior probability is **proportion times** the prior

The posterior probability is proportional to the numerator, i.e., the likelihood



2. Can you describe the difference between discrete and continuous variables?

discrete variables \rightarrow have a set of predefined values, an integer e.g., having a tail vs. not

continuous variables \rightarrow can take on any real number value within a range e.g., length, body mass





2. Can you describe the difference between probabilities and probability densities?

probabilities \rightarrow a probability takes a singular value, e.g. P = 0.5

probability densities \rightarrow a range of values represented by a distribution



3. What is the difference between vague vs. informative priors?

value is e.g., it could be anything between 0 and infinity

 \rightarrow an informative prior is used for parameters where we have some good existing knowledge about what the parameter value could be prior distribution with a mean equal to the known value and add a small variance

 \rightarrow a vague prior is used for parameters where we have little clue what the true

e.g., maybe we already know the rate of evolution among a well studied group, so we could use a



4. What is the aim of MCMC in Bayesian inference?

→ the aim is to **approximate** the posterior distribution

use MCMC to traverse the parameter space and at each step calculate the areas with the highest posterior probability

The posterior distribution is hard to calculate analytically (i.e., exactly), so we likelihood x the prior, spending time in different regions of the parameter space in proportion to their posterior probability - this means, we spend most time in



Bayesian tree inference

posterior





Bayesian tree inference

$\mathsf{P}\left(\begin{smallmatrix} 0101...\\1101...\\0100...\end{smallmatrix}\right) \mathsf{P}\left(\begin{smallmatrix} \mathbf{P} \\ \mathbf{P} \\$



this part is incredibly difficult to calculate!



Hastings ratio

new parameter values









What is Markov chain Monte Carlo (MCMC)?



The aim is to produce a histogram that provides a good approximation of the posterior



Q&A Phylogenetics primer part 3b by Paul Lewis

(the answers provided here are my in may vary!)

(the answers provided here are my interpretation of these concepts – answers



1. How are steps chosen in an MCMC analysis?

→ these depend on the type of parameter and the landscape of the parameters space





2. Give an example of a parameter you would estimate under each of the following prior distributions and try to state why?

Gamma distribution

Lognormal distribution

Beta distribution

Dirichlet distribution



3. Why do we sometimes need to calculate the marginal likelihood?

→ this is required for model testing within a Bayesian framework



4. What is the difference between a hierarchical model and a non-hierarchical model?

 \rightarrow in a hierarchical model different components of the model are nested, that apply to the data

e.g., different models can be joined together to model different processes

