

Phylogenetic Inference using RevBayes

Discrete Morphology

April M. Wright and Michael J. Landis

Introduction

While molecular data have become the default for building phylogenetic trees for many types of evolutionary analysis, morphological data remains important, particularly for analyses involving fossils. The use of morphological data raises special considerations for model-based methods for phylogenetic inference. Morphological data are typically collected to maximize the number of parsimony-informative characters - that is, the characters that favor one topology over another. Morphological characters also do not carry common meanings from one character in a matrix to the next; character codings are made arbitrarily. These two factors require extensions to our existing phylogenetic models. Accounting for the complexity of morphological characters remains challenging. This tutorial will provide a discussion of modeling morphological characters, and will demonstrate how to perform Bayesian phylogenetic analysis with morphology using **RevBayes**.

Contents

The Discrete Morphology guide contains several tutorials

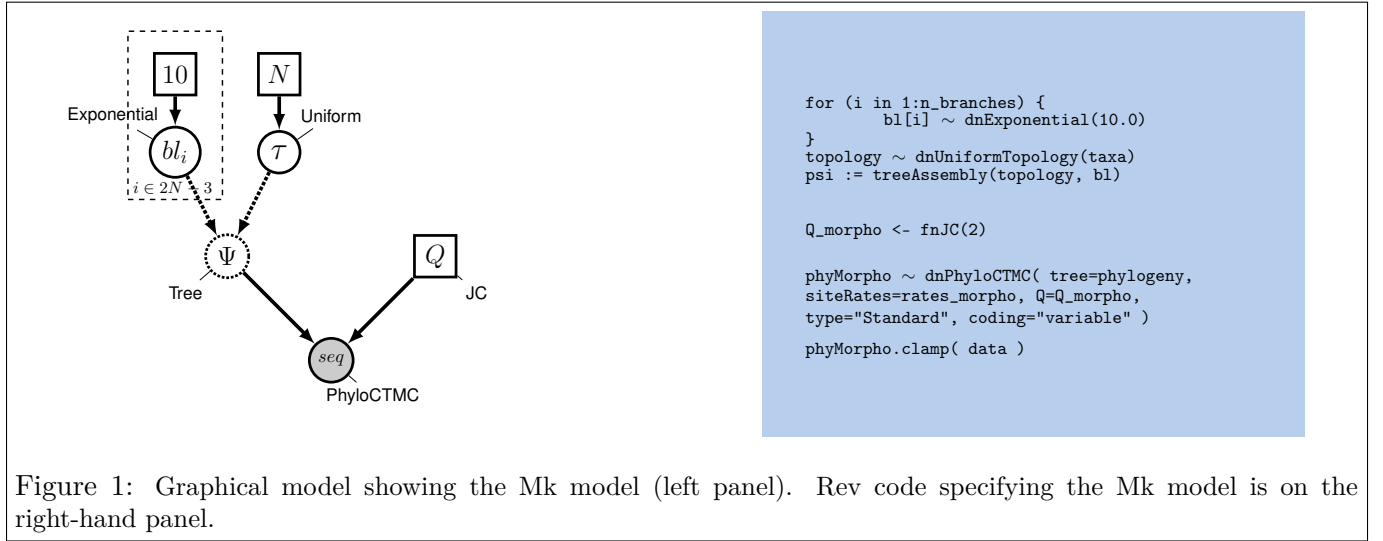
- Section 1: Overview of the Discrete Morphological models
- Section 2: A simple discrete morphology analysis
- Section 3: A model for allowing state frequency variation across binary characters
- Section 4: A model for allowing state frequency variation across binary and multistate characters
- Section 5: Evaluating the MCMC

Recommended tutorials

The Discrete Morphology tutorials assume the reader is familiar with the content covered in the following RevBayes tutorials

- Rev Basics
- Molecular Models of Character Evolution
- Running and Diagnosing an MCMC Analysis
- Divergence Time Estimation and Node Calibrations

1 Overview of Discrete Morphology Models



Molecular data forms the basis of most phylogenetic analyses today. However, morphological characters remain relevant: Fossils often provide our only direct observation of extinct biodiversity; DNA degradation can make it difficult or impossible to obtain sufficient molecular data from fragile museum specimens. Using morphological data can help researchers include specimens in their phylogeny that might be left out of a molecular tree.

To understand how morphological characters are modeled, it is important to understand how characters are collected. Unlike in molecular data, for which homology is algorithmically determined, homology in a character is typically assessed an expert. Biologists will typically decide what characters are homologous by looking across specimens at the same structure in multiple taxa; they may also look at the developmental origin of structures in making this assessment (Phillips 2006). Once homology is determined, characters are broken down into states, or different forms a single character can take. The state ‘0’ commonly refers to absence, meaning that character is not present. In some codings, absence will mean that character has not evolved in that group. In others, absence means that that character has not evolved in that group, and/or that that character has been lost in that group (Freudenstein 2005). This type of coding is arbitrary, but both **non-random** and **meaningful**, and poses challenges for how we model the data.

Historically, most phylogenetic analyses using morphological characters have been performed using the maximum parsimony optimality criterion. Maximum parsimony analysis involves proposing trees from the morphological data. Each tree is evaluated according to how many changes it implied in the data, and the tree that requires the fewest changes is preferred. In this way of estimating a tree, a character that does not change, or changes only in one taxon, cannot be used to discriminate between trees (i.e., it does not favor a topology). Therefore, workers with parsimony typically do not collect characters that are parsimony uninformative.

In 2001, Paul Lewis (Lewis 2001) introduced a generalization of the Jukes-Cantor model of sequence evolution for use with morphological data. This model, called the Mk (Markov model, assuming each character is in one of k states) model provided a mathematical formulation that could be used to estimate trees from morphological data in both likelihood and Bayesian frameworks. While this model is a useful step forward, as a generalization of the Jukes-Cantor, it still makes fairly simplistic assumptions. This tutorial will guide you through estimating a phylogeny with the Mk model, and two useful extensions to

the model.

1.1 The Mk Model

Make a copy of the MCMC and model files you just made. Call them `mcmc_mk_discretized.Rev` and `model_mk_discretized.Rev`. These will contain the new model parameters and models.

The Mk model is a generalization of the Jukes-Cantor model of nucleotide sequence evolution, which we discussed in **Molecular Models of Character Evolution**. The Q matrix for a two-state Mk model looks like so:

$$Q = \begin{pmatrix} -\mu_0 & \mu_{01} \\ \mu_{10} & -\mu_1 \end{pmatrix},$$

This matrix can be expanded to accommodate multi-state data, as well:

$$Q = \begin{pmatrix} -\mu_0 & \mu_{01} & \mu_{02} & \mu_{03} \\ \mu_{10} & -\mu_1 & \mu_{12} & \mu_{13} \\ \mu_{20} & \mu_{21} & -\mu_2 & \mu_{23} \\ \mu_{30} & \mu_{31} & \mu_{32} & -\mu_3 \end{pmatrix},$$

However, the Mk model sets transitions to be equal from any state to any other state. In that sense, our multistate matrix really looks like this:

$$Q = \begin{pmatrix} -\mu_0 & \mu & \mu & \mu \\ \mu & -\mu_1 & \mu & \mu \\ \mu & \mu & -\mu_2 & \mu \\ \mu & \mu & \mu & -\mu_3 \end{pmatrix},$$

Because this is a Jukes-Cantor-like model, state frequencies do not vary as a model parameter. These assumptions may seem unrealistic. However, all models are a compromise between reality and generalizability. Prior work has demonstrated that, in many conditions, the model does perform adequately (Wright and Hillis 2014). Because morphological characters do not carry common meaning across sites in a matrix in the way that nucleotide characters do, making assumptions that fit all characters is challenging. A visualization of this simple model can be seen in Fig. 1.

We will first perform a phylogenetic analysis using the Mk model. In further sections, we will explore how to relax key assumptions of the Mk model.

1.2 Ascertainment Bias

When Lewis first introduced the Mk model, he observed that branch lengths on the trees were greatly inflated. The reason for this is that when morphological characters are collected, characters that do not vary, or vary in a non-parsimony-informative way (such as autapomorphies) are excluded. Excluding these

low-rate characters causes the overall amount of evolution to be over-estimated. This causes an inflation in the branch lengths [Lewis \(2001\)](#).

Therefore, when performing a morphological phylogenetic analysis, it is important to correct for this bias. There are numerous statistically valid ways to perform this correction [Allman and Rhodes \(2008\)](#). Original corrections simulated invariant and non-parsimony informative characters along the proposed tree. The likelihood of these characters would then be calculated and used to normalize the total likelihood value. RevBayes implements a dynamic programming approach that calculates the same likelihood, but does so faster.

2 Example: Inferring a Phylogeny of Fossil Bears Using the Mk Model

In this example, we will use morphological character data from 18 taxa of extinct bears ([Abella et al. 2011](#)). The dataset contains 62 binary characters, a fairly typical dataset size for morphological characters.

2.1 Tutorial Format

This tutorial follows a specific format for issuing instructions and information.

The boxed instructions guide you to complete tasks that are not part of the RevBayes syntax, but rather direct you to create directories or files or similar.

Information describing the commands and instructions will be written in paragraph-form before or after they are issued.

All command-line text, including all Rev syntax, are given in **monotype font**. Furthermore, blocks of Rev code that are needed to build the model, specify the analysis, or execute the run are given in separate shaded boxes. For example, we will instruct you to create a constant node called **example** that is equal to 1.0 using the `<-` operator like this:

```
example <- 1.0
```

It is important to be aware that some PDF viewers may render some characters given as **Rev commands** differently. Thus, if you copy and paste text from this PDF, you may introduce some incorrect characters. Because of this, we recommend that you type the instructions in this tutorial or copy them from the scripts provided.

2.2 Data and Files

On your own computer, create a directory called **RB_DiscreteMorphology_Tutorial** (or any name you like).

In this directory download and unzip the archive containing the data files: [data.zip](#).

This will create a folder called **data** that contains the files necessary to complete this exercise.

2.3 Getting Started

Create a new directory (in **RB_DiscreteMorphology_Tutorial**) called **scripts**. (If you do not have this folder, please refer to the directions in section 2.2.)

When you execute **RevBayes** in this exercise, you will do so within the main directory you created (**RB_DiscreteMorphology_Tutorial**), thus, if you are using a Unix-based operating system, we recommend that you add the **RevBayes** binary to your path.

2.4 Creating Rev Files

For complex models and analyses, it is best to create **Rev** script files that will contain all of the model parameters, moves, and functions. In this exercise, you will work primarily in your text editor¹ and create a set of modular files that will be easily managed and interchanged. In this first section, you will write the following files from scratch and save them in the **scripts** directory:

- **mcmc_mk.Rev**: the master **Rev** file that loads the data, the separate model files, and specifies the monitors and MCMC sampler.
- **model_mk.Rev**: specifies the model describing discrete morphological character change (binary characters).

All of the files that you will create are also provided in the **RevBayes** tutorial repository². Please refer to these files to verify or troubleshoot your own scripts.

Open your text editor and create the master **Rev** file called **mcmc_Mk.Rev** in the **scripts** directory.

Enter the **Rev** code provided in this section in the new model file.

The file you will begin in this section will be the one you load into **RevBayes** when you've completed all of the components of the analysis. In this section you will begin the file and write the **Rev** commands for loading in the taxon list and managing the data matrices. Then, starting in section ??, you will move on to writing module files for each of the model components. Once the model files are complete, you will return to editing **mcmc_Mk.Rev** and complete the **Rev** script with the instructions given in section 2.6.

2.4.1 Load Data Matrices

RevBayes uses the function **readDiscreteCharacterData()** to load a data matrix to the workspace from a formatted file. This function can be used for both molecular sequences and discrete morphological characters. Import the morphological character matrix and assign it to the variable **morpho**.

¹In section ?? we offer a recommendation for a text editor.

²https://github.com/revbayes/revbayes_tutorial/tree/master/RB_Discrete_Morphology_Tutorial/scripts

```
morpho <- readDiscreteCharacterData("data/bears.nex")
```

2.4.2 Create Helper Variables

Before we begin writing the **Rev** scripts for each of the model components, we need to instantiate a couple “helper variables” that will be used by downstream parts of our model specification files. These variables will be used in more than one of the module files so it’s best to initialize them in the master file.

Create a new constant node called **n_taxa** that is equal to the number of species in our analysis (18). We will also create a constant node of the taxon names. This list will be used to initialize the tree.

```
taxa <- morpho.names()
n_taxa <- morpho.size()
```

Next, create a workspace variable called **mvi**. This variable is an iterator that will build a vector containing all of the MCMC moves used to propose new states for every stochastic node in the model graph. Each time a new move is added to the vector, **mvi** will be incremented by a value of 1.

```
mvi = 1
```

One important distinction here is that **mvi** is part of the **RevBayes** workspace and not the hierarchical model. Thus, we use the workspace assignment operator **=** instead of the constant node assignment **<-**.

Save your current working version of **mcmc_Mk.Rev** in the **scripts** directory.

We will now move on to the next **Rev** file and will complete **mcmc_Mk.Rev** in section 2.6.

2.5 The Mk Model

Open your text editor and create the master **Rev** file called **model_Mk.Rev** in the **scripts** directory.

Enter the **Rev** code provided in this section in the new model file.

First, we will create a vector of moves on branch lengths. This should be familiar from the **RB_CTMC** tutorial:

```

nbr <- 2*names.size() - 3
for (i in 1:nbr){
  br_lens[i] ~ dnExponential(5.0)
  moves[mvi++] = mvScale(br_lens[i])
}

```

Next, we will create a Q matrix. Recall that the Mk model is simply a generalization of the JC model. Therefore, we will create a 2x2 Q matrix using **fnJC**, which initializes Q-matrices with equal transition probabilities between all states.

```

Q_morpho := fnJC(2)

```

Now that we have the basics of the model specified, we will add Gamma-distributed rate variation and specify moves on the parameter to the Gamma distribution.

```

alpha_morpho ~ dnExponential( 1.0 )
rates_morpho := fnDiscretizeGamma( alpha_morpho, alpha_morpho, 4 )

#Moves on the parameters to the Gamma distribution.
moves[mvi++] = mvScale(alpha_morpho, lambda=0.01, weight=1.0)
moves[mvi++] = mvScale(alpha_morpho, lambda=0.1, weight=1.0)
moves[mvi++] = mvScale(alpha_morpho, lambda=1, weight=1.0)

```

Next we assemble the tree and specify a move on the topology.

```

tau ~ dnUniformTopology(names)
phylogeny := treeAssembly(tau, br_lens)
moves[mvi++] = mvSPR(tau)

```

Lastly, we set up the CTMC. This should be familiar from the **RB_CTMC** tutorial. We see some familiar pieces: tree, Q matrix and site_rates. We also have two new keywords: data type and coding. The data type argument specifies the type of data - in our case, "Standard", the specification for morphology. Coding specifies what type of ascertainment bias is expected. We are using the 'variable' correction, as we have no invariant character in our matrix. If we also lacked parsimony non-informative characters, we would use the coding 'informative'.

```

phyMorpho ~ dnPhyloCTMC(tree=phylogeny, siteRates=rates_morpho, Q=Q_morpho, type="
  Standard", coding="variable")
phyMorpho.clamp(morpho)

```

All of the components of the model are now specified.

2.6 Complete Master Rev File

Return to the master Rev file you created in section ?? called **mcmc_Mk.Rev** in the **scripts** directory.

Enter the Rev code provided in this section in this file.

2.6.1 Source Model Scripts

RevBayes uses the **source()** function to load commands from Rev files into the workspace. Use this function to load in the model scripts we have written in the text editor and saved in the **scripts** directory.

```
source("scripts/model_Mk.Rev")
```

2.6.2 Create Model Object

We can now create our workspace model variable with our fully specified model DAG. We will do this with the **model()** function and provide a single node in the graph (**phylogeny**).

```
mymodel = model(phylogeny)
```

The object **mymodel** is a wrapper around the entire model graph and allows us to pass the model to various functions that are specific to our MCMC analysis.

2.6.3 Specify Monitors and Output Filenames

The next important step for our master Rev file is to specify the monitors and output file names. For this, we create a vector called **monitors** that will each sample and record or output our MCMC.

First, we will specify a workspace variable to iterate over the **monitors** vector.

```
mni = 1
```

The first monitor we will create will monitor every named random variable in our model graph. This will include every stochastic and deterministic node using the **mnModel** monitor. The only parameter that is not included in the **mnModel** is the tree topology. Therefore, the parameters in the file written by this monitor are all numerical parameters written to a tab-separated text file that can be opened by accessory programs for evaluating such parameters. We will also name the output file for this monitor and indicate that we wish to sample our MCMC every 10 cycles.


```
monitors[mni++] = mnModel(filename="output/mk_simple.log", printgen=10)
```

The **mnFile** monitor writes any parameter we specify to file. Thus, if we only cared about the branch lengths and nothing else (this is not a typical or recommended attitude for an analysis this complex) we wouldn't use the **mnModel** monitor above and just use the **mnFile** monitor to write a smaller and simpler output file. Since the tree topology is not included in the **mnModel** monitor (because it is not numerical), we will use **mnFile** to write the tree to file by specifying our **phylogeny** variable in the arguments.

```
monitors[mni++] = mnFile(filename="output/mk_simple.trees", printgen=10, phylogeny)
```

The last monitor we will add to our analysis will print information to the screen. Like with **mnFile** we must tell **mnScreen** which parameters we'd like to see updated on the screen.

```
monitors[mni++] = mnScreen(printgen=10)
```

2.6.4 Set-Up the MCMC

Once we have set up our model, moves, and monitors, we can now create the workspace variable that defines our MCMC run. We do this using the **mcmc()** function that simply takes the three main analysis components as arguments.

```
mymcmc = mcmc(myModel, monitors, moves)
```

The MCMC object that we named **mymcmc** has a member method called **.run()**. This will execute our analysis and we will set the chain length to **10000** cycles using the **generations** option.

```
mymcmc.run(generations=100000)
```

Once our Markov chain has terminated, we will want RevBayes to close. Tell the program to quit using the **q()** function.

```
q()
```

You made it! Save all of your files.

2.7 Execute the MCMC Analysis

With all the parameters specified and all analysis components in place, you are now ready to run your analysis. The `Rev` scripts you just created will all be used by `RevBayes` and loaded in the appropriate order.

Begin by running the `RevBayes` executable. In Unix systems, type the following in your terminal (if the `RevBayes` binary is in your path):

```
rb
```

Provided that you started `RevBayes` from the correct directory (`RB_DiscreteMorphology_Tutorial`), you can then use the `source()` function to feed `RevBayes` your master script file (`mcmc_mk.Rev`).

```
source("scripts/mcmc_mk.Rev")
```

This will execute the analysis and you should see the following output (though not the exact same values):

```
source("scripts/mcmc_mk.rev")
  Processing file "scripts/mcmc_mk.rev"
  Successfully read one character matrix from
  file 'data/bears.nex'
  Processing file "scripts/mk_simple.rev"
  Processing of file "scripts/mk_simple.rev"
  completed

  Running MCMC simulation
  This simulation runs 1 independent replicate.
  The simulator uses 37 different moves in a
  random move schedule with 37 moves per iteration
```

Iter	Prior	Posterior elapsed	Likelihood ETA
0	-685.779	-666.105	-19.6731
10	-633.737	-608.951	-24.786
20	-593.634	-558.797	-34.8368
30	-578.661	-536.125	-42.5362
40	-544.22	-514.098	-30.1223
50	-515.505	-492.988	-22.5169
60	-483.695	-461.347	-22.3479
...			

When the analysis is complete, `RevBayes` will quit and you will have a new directory called `output` that will contain all of the files you specified with the monitors (Sect. 2.6.3).

3 Example: Relaxing the Assumption of Equal Transition Probabilities

Make a copy of the MCMC and model files you just made. Call them `mcmc_mk_discretized.Rev` and `model_mk_discretized.Rev`. These will contain the new model parameters and models.

Figure 2: Graphical model demonstrating the discretized Beta distribution for allowing variable state frequencies.

The Mk model makes a number of assumptions, but one that may strike you as unrealistic is the assumption that characters are equally likely to change from any one state to any other state. That means that a trait is as likely to be gained as lost. While this may hold true for some traits, we expect that it may be untrue for many others.

RevBayes has functionality to allow us to relax this assumption. We do this by specifying a Beta prior on state frequencies. Remember from the **RB_CTMC** lesson that stationary frequencies impact how likely we are to see changes in a character. For example, it may be very likely, in a character, to change from 0 to 1. But if the frequency of 0 is very low, we will still seldom see this change.

We can exploit the relationship between state frequencies and observed changes to allow for variable Q matrices across characters (Fig. 2). To do this, we generate a Beta distribution on state frequencies ([Huelsenbeck and Ronquist 2001](#)), and use the state frequencies from that Beta distribution to generate a series of Q-matrices to use to evaluate our data.

This type of model is called a **mixture model**. There are assumed to be subdivisions in the data, which may require different parameters (in this case, state frequencies). These subdivisions are not defined *a priori*. This model has previously been shown to be effective for a range of empirical and simulated datasets ([Wright et al. 2016](#)).

3.1 Modifying the MCMC File

At each place in which the output files are specified in the MCMC file, change the output path so you don't overwrite the output from the previous exercise. For example, you might call your output file `output/mk_discretized.log` and `output/mk_discretized.trees`. Change source statement to indicate the new model file.

3.2 Modifying the Model File

Open the new model file that you created. We need to modify the way in which the Q matrix is specified. We will use a discretized Beta distribution to place a prior on state frequencies. The Beta distribution has two parameters, α and β . These two parameters specify the shape of the distribution. State frequencies will be evaluated according to this distribution, in the same way that rate variation is evaluated according to the Gamma distribution. The discretized distribution is split into multiple classes, each with its own set of frequencies for the 0 and 1 characters. The number of classes can vary; we have chosen 4 for tractability.

```
n_cats = 4
alpha_ofbeta ~ dnExponential( 1 )
beta_ofbeta ~ dnExponential( 1 )
moves[mvi++] = mvScale(alpha_ofbeta, lambda=1, weight=1.0 )
moves[mvi++] = mvScale(alpha_ofbeta, lambda=0.1, weight=1.0 )
```

```
moves[mvi++] = mvScale(alpha_ofbeta, lambda=0.01, weight=1.0 )
moves[mvi++] = mvScale(beta_ofbeta, lambda=1, weight=1.0 )
moves[mvi++] = mvScale(beta_ofbeta, lambda=0.1, weight=1.0 )
moves[mvi++] = mvScale(beta_ofbeta, lambda=0.01, weight=1.0 )
```

Above, we initialized the number of categories, the parameters to the Beta distribution, and the moves on the parameters to the Beta.

Next, we set the categories to each represent a quadrant of the Beta distribution specified by the **alpha_ofbeta** and **beta_ofbeta**.

```
cats := fnDiscretizeBeta(alpha_ofbeta, beta_ofbeta, 4)
```

If you were to print the **cats** variable, you would see a list of state frequencies like so:

```
0.0780943
0.40457
0.769453
0.97106
```

Using these state frequencies, we will generate a new vector of **Q** matrices. Because we are varying the state frequencies, we must use a **Q** matrix generation function that allows for state frequencies to vary as a parameter. We will, therefore, use the **freeBinary** function.

```
for (i in 1:cats.size())
{
Q[i] := fnFreeBinary(v(1-cats[i], cats[i]), false)
}
```

Once we've made the our vector of matrices, we specified moves on our matrix vector:

```
matrix_probs ~ dnDirichlet(v(1,1,1,1))
moves[mvi++] = mvSimplexElementScale(matrix_probs, alpha=10, weight=1.0)
```

This Dirichlet prior says that no category is expected to have more characters than another. If you expected some category to hold more of the characters, you could put more weight on that category.

The only other specification that needs to change in the model file is the CTMC:

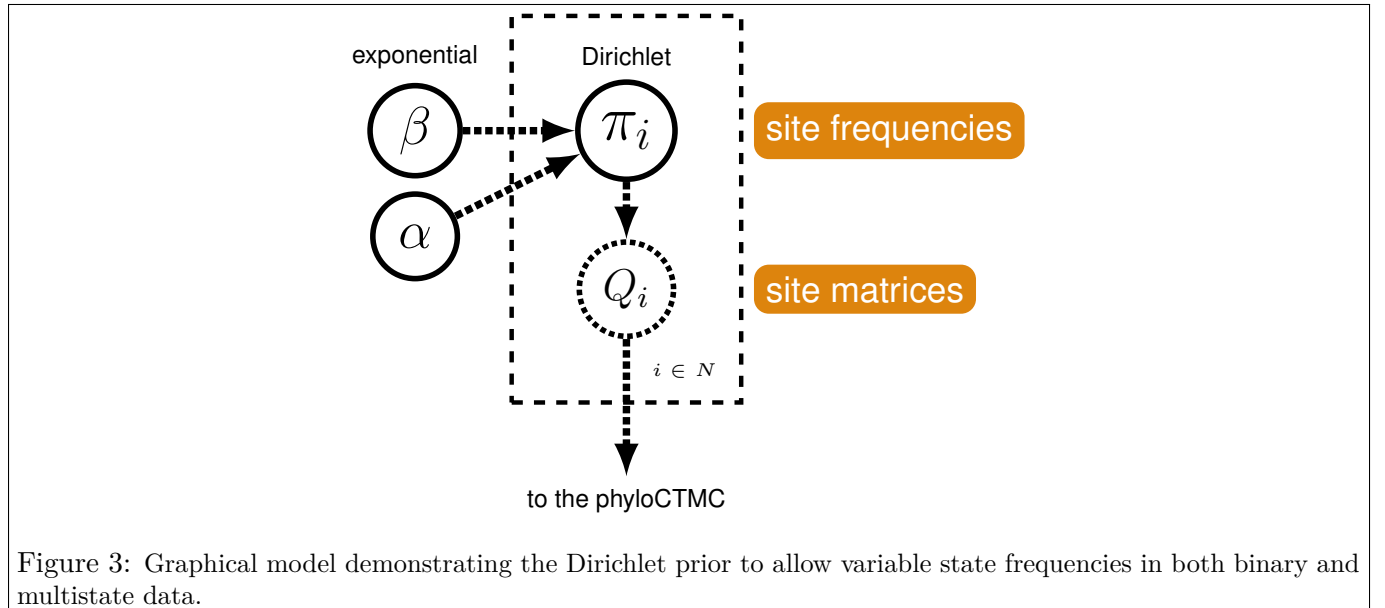
```
phyMorpho ~ dnPhyloCTMC(tree=phylogeny, siteRates=rates_morpho, Q=Q, type="Standard",
    coding="variable", siteMatrices=matrix_probs)
```

You'll notice that we've added a command to tell the CTMC that we have multiple site matrices that will be applied to different characters in the matrix.

3.2.1 Set-Up the MCMC

The MCMC chain set-up does not need to change. Run the new MCMC file, just as you ran the plain Mk file. This estimation will take longer than the Mk model, due to increased model complexity.

4 Site-Heterogeneous Discrete Morphology Model



In the previous example, we explored allowing among-character variation in state frequencies. This is an excellent start for allowing more complex models for morphology. But this approach also has several shortcomings. First, because we use a Beta distribution, this model really only works for binary data. Secondly, oftentimes, we will not have a good idea of the shape of the distribution from which we expect state frequencies to be drawn.

To accommodate for these concerns, **RevBayes** also has a model that is similar to the CAT model (Lartillot and Philippe 2004).

The site-heterogeneous discrete morphology model (SHDM) uses a hyperprior on the prior on state frequencies to mix over different possible combinations state frequencies. In this mixture model, F81 Q-matrices (an extension of the Jukes-Cantor which allows for different state frequencies between characters) is initialized from a set of state frequencies. The number of Q-matrices initialized is equal to the number of user-defined categories, as in the discretized Beta model. The state frequencies used to initialize the Q-matrices are drawn from a Dirichlet prior distribution, which is generated by drawing values from an exponential hyperprior distribution. This model is visualized in Fig. 3.

4.1 Example: Site-Heterogeneous Discrete Morphology Model

Make a copy of the MCMC and model files you just made. Call them `mcmc_mk_hyperprior.Rev` and `model_mk_hyperprior.Rev`. These will contain the new model parameters and models.

4.2 Modifying the MCMC File

At each place in which the output files are specified in the MCMC file, change the output path so you don't overwrite the output from the previous exercise. For example, you might call your output file `output/mk_hyperprior.log` and `output/mk_hyperprior.trees`. We will also monitor `Q_morpho` and `pi`. Add `Q_morpho` and `pi` to the `mnScreen`. Change source statement to indicate the new model file.

4.3 Modifying the Model File

Open the new model file that you created. We need to modify the way in which the Q-matrix is specified. First, we will create a hyperprior called `dir_alpha` and specify a move on it.

```
dir_alpha ~ dnExponential(1)
moves[mvi++] = mvScale(dir_alpha, lambda=1, weight=1.0 )
moves[mvi++] = mvScale(dir_alpha, lambda=0.1, weight=1.0 )
moves[mvi++] = mvScale(dir_alpha, lambda=0.01, weight=1.0 )
```

This hyperparameter, `dir_alpha`, will be used as a parameter to a Dirichlet distribution from which our state frequencies will be drawn.

```
pi_prior := v(dir_alpha, dir_alpha)
```

If you were using multistate data, the `dir_alpha` can be repeated for each state. Next, we will modify our previous loop to use these state frequencies to initialize our Q-matrices.

```
{
    pi[i] ~ dnDirichlet(pi_prior)
    moves[mvi++] = mvSimplexElementScale(pi[i], alpha=10, weight=1.0)

    Q_morpho[i] := fnF81(pi[i])
}
```

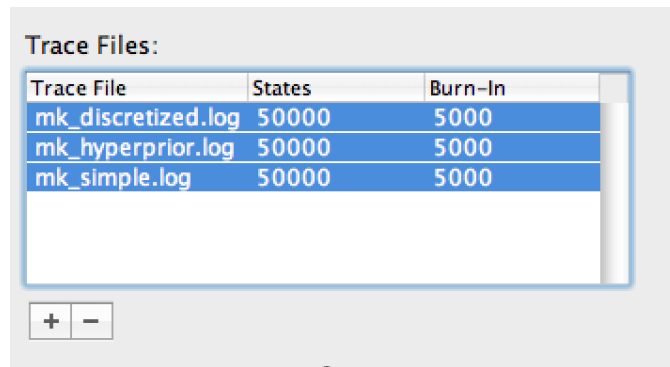
In the above loop, for each of our categories, we make a new draw of state frequencies from our Dirichlet distribution (the shape of which is determined by our `dir_alpha` values). We then use `fnF81` to make our Q-matrices. For each **RevBayes** iteration, we will have 4 pi values and 4 Q-matrices, one for each of the number of categories we specified.

No other aspects of the model file need to change. Run the MCMC as before.

5 Evaluate and Summarize Your Results

5.1 Evaluate MCMC

As discussed in **RB_Total_Evidence_TEFBD**, we will use Tracer to evaluate the MCMC samples from our three estimations. Load all three of the MCMC logs into the Tracer window. The MCMC chains will not have converged because they have not been run very long. Highlight all three files in the upper left-hand viewer (Fig. 4) by right- or command-clicking all three files.



The screenshot shows the 'Trace Files' window in Tracer. It contains a table with three columns: 'Trace File', 'States', and 'Burn-In'. Three rows are listed, all of which are highlighted in blue: 'mk_discretized.log', 'mk_hyperprior.log', and 'mk_simple.log'. Each row shows '50000' for 'States' and '5000' for 'Burn-In'. Below the table are '+' and '-' buttons.

Trace File	States	Burn-In
mk_discretized.log	50000	5000
mk_hyperprior.log	50000	5000
mk_simple.log	50000	5000

Figure 4: Highlight all three files for model comparison.

Once all three trace logs are loaded and highlighted, first look at the estimated marginal likelihoods. You will notice that the Mk model, as originally proposed by Lewis (2001) is improved by allowing any state frequency heterogeneity at all. The discretized model and the Dirichlet model both represent improvements, but are fairly close in likelihood score to each other (Fig. 5). Likely, we would need to perform stepping stone model assessment to truly tell if the more complicated model is statistically justified. This analysis is too complicated and time-consuming for this tutorial period, but you will find instructions below for performing the analysis.

Click on the ‘Trace’ panel. In the lower left hand corner, you will notice an option to color each trace by the file it came from. Choose this option (you may need to expand the window slightly to see it). Next to this option, you can also see an option to add a legend to your trace window. The results of this coloring can be seen in Fig. 6. When the coloring is working, you will see that the Mk model mixes quite well, but that mixing becomes worse as we relax the assumption of equal state frequencies. This is because we are greatly increasing model complexity. Therefore, we would need to run the MCMC chains longer if we were to use these analyses in a paper.

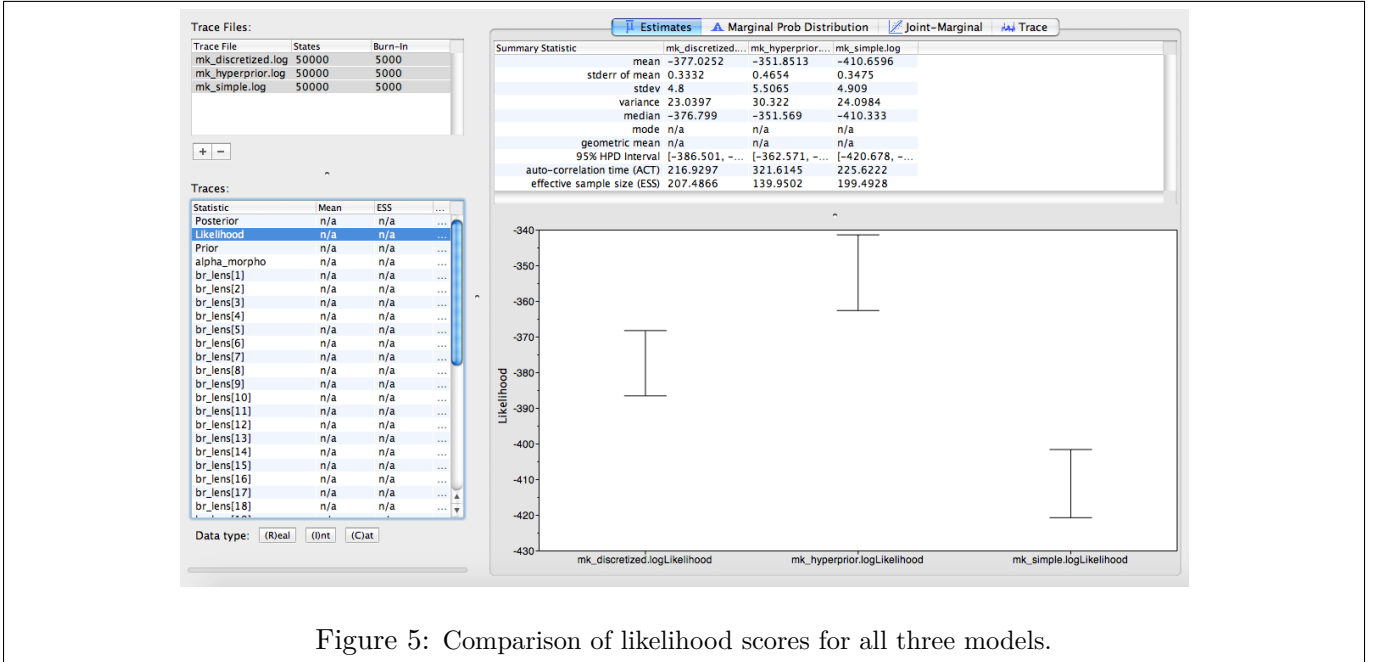


Figure 5: Comparison of likelihood scores for all three models.

Figure 6: The Trace window. The traces are colored by which version of the Mk model they correspond to.

6 Ancestral state estimation

Discrete morphological models are not only useful for tree estimation, as was done in Section 1, but also for ancestral state estimation. The central problem in statistical phylogenetics concerns *marginalizing* over all unobserved character histories that evolved along the branches of a given phylogenetic tree. This allows one to compute the likelihood of observing a character matrix given they evolved according to an underlying tree. While marginalizing over those histories, one might also seek to produce samples of ancestral states that are consistent with the tip state distribution.

$$P(X_{\text{tip}}) = \sum_{x \in X_{\text{internal}}} P(X_{\text{tip}}, X_{\text{internal}})$$

We want the joint probabilities of internal node states conditioned on the tips

$$P(X^{(\text{internal})} | X^{(\text{tip})}) = \sum_{x \in X_{\text{internal}}} P(X_{\text{tip}}, X_{\text{internal}})$$

or the marginal probabilities conditioned on the tips

$$P(X_i^{(\text{internal})} | X^{(\text{tip})}) = \sum_{x \in X_{\text{internal}}} P(X_{\text{tip}}, X_{\text{internal}})$$

Ancestral state estimation is the procedure of sampling ancestral states that are consistent with the observed distribution of tip states. There are several strategies for doing so, each with their own advantages and disadvantages. One can compute the *exact probability* of any given ancestral node having an ancestral state value; or one can *sample* states from that distribution in proportion to the exact probability. One can produce *marginal* ancestral state estimates that report how often any given node is in a state without the context of the ancestral states of the remaining nodes; or one can produce *joint* ancestral state estimates that report the probability of particular combinations of sequences of evolutionary histories, running from the root to the tips of the phylogeny. While it is straightforward to produce exact probabilities for marginal histories, sampling is generally required to produce joint histories.

6.1 Ancestral state monitors

We focus on joint samples in this tutorial and in `RevBayes` because they yield most direct representation of possible evolutionary scenarios. The completed exercises from Section 1 made use of a monitor names `mnJointConditionalAncestralState`, which was created with the command

[illegible]

The core arguments this monitor needs are a tree object (`tree=phylogeny`), the phylogenetic model (`ctmc=phyMorpho`), an output filename (`filename="output/simple.states.txt"`), the data type for the characters (`type="Standard"`), and the sampling frequency (`printgen=10`). The final argument, `withTips=false`, indicates that we do not wish to record the starting state for each branch, since it will always be identical to the ending state of the parental branch (but this is not necessarily so for models with cladogenesis).

The monitor will produce a joint sample of ancestral states (conditional on the tip states) every 10 iterations, storing the samples to the file "output/simple.states.txt". Viewing the states file, we see

Iteration	end_1	end_2	end_3	end_4
0	? , 0 , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , 0 , 0 , 0 , ...			
10	? , 0 , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , 0 , 0 , 0 , ...			
20	? , 0 , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , 0 , 0 , 0 , ...			
30	? , 0 , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , 0 , 0 , 0 , ...			
40	? , 0 , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , 0 , 0 , 0 , ...			
50	? , 0 , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , ? , 0 , 0 , 0 , ...			
...				

The first column is the MCMC iteration, and the remaining columns represent the states sampled at end of the branch for an indexed node. When the phylogeny is a random variable, the index is used to map an ancestral state estimate to the node in any given tree sample from the analysis. The vector of comma-delimited ancestral state estimates is reported for each column value. In the example above, we see that

the node indexed 1 (a tip node) always produces the observed character matrix [future versions will sample ambiguous tip states, e.g. ?]. The first $1 \dots N$ index values are assigned to tip nodes in a tree with N tips. We observe the uncertainty in our ancestral state estimates by viewing the final column

```
[ column headers are far to the left ]
... 1,0,1,0,1,0,1,1,0,0,0,1,1,1
... 1,1,1,1,1,0,0,1,0,0,1,1,1,0
... 1,1,1,0,0,0,0,0,0,0,0,0,1,0
... 1,1,0,1,1,0,1,0,0,1,0,0,1,0
... 1,1,0,1,1,0,1,0,1,0,0,0,1,0
... 1,1,0,1,1,0,1,0,1,1,0,0,1,0
...
```

6.2 Ancestral state figures

Now that we have our posterior distribution of ancestral states and of phylogenetic trees, we want to summarize those results into something simpler. This section will aim to produce a pdf containing figures for the ancestral state estimates mapped onto a consensus tree for all 62 characters. To accomplish this, we will rely on a combination of RevBayes, R, and bash scripts. [Unfortunately, this exact pipeline probably does not work for Windows users at this time.]

To generate the figure `output/mksimple`.

```
./scripts/make_anc_trees.sh output/mk_simple 1 5
```

7 Beyond binary rate matrices

The instantaneous rate matrix encodes the transition rates between all pairs of evolutionary states. It is important to emphasize that all rate matrices are assertions about how morphological evolution operates. Depending on how one populates the rate matrix elements, different evolutionary hypotheses may be expressed.

7.1 Symmetric unordered

The standard Mk model of character evolution, where M denotes it is a Markov model and K denotes the number of states for the character. The lineage may transition directly from state 1 to state 4 without going through states 2 and 3, which is representative of a character with *unordered* states. In addition, all transition rates are equal as they are in the Jukes-Cantor rate matrix (Jukes and Cantor 1969). Here is an example of a symmetric unordered Mk model for $K = 4$.

$$Q = \begin{pmatrix} - & r & r & r \\ r & - & r & r \\ r & r & - & r \\ r & r & r & - \end{pmatrix}$$

Define the single shared rate parameter

```
r <- 1.0
```

Define the rates

```
rates := [ [0.0, r,  r,  r],
           [ r, 0.0, r,  r],
           [ r,  r, 0.0, r],
           [ r,  r,  r, 0.0] ]
```

Create the rate matrix

```
Q := fnFreeK(rates)
Q
[ [ -1.0000, 0.3333, 0.3333, 0.3333 ] ,
  [ 0.3333, -1.0000, 0.3333, 0.3333 ] ,
  [ 0.3333, 0.3333, -1.0000, 0.3333 ] ,
  [ 0.3333, 0.3333, 0.3333, -1.0000 ] ]
```

Compute the transition probability matrix for a branch length of 0.1.

```
P <- Q.getTransitionProbabilities(0.1)
P
[ [ 0.906, 0.031, 0.031, 0.031],
  [ 0.031, 0.906, 0.031, 0.031],
  [ 0.031, 0.031, 0.906, 0.031],
  [ 0.031, 0.031, 0.031, 0.906] ]
```

7.2 Asymmetric ordered

Character states that are ordered imply that evolutionary transitions occur in particular sequences. For example, the number of digits on a foot might vary by gaining and losing single digits, meaning the transition from three to five digits cannot occur without going through the evolutionary state of possessing four digits. Below, we assume that gain events ($n \rightarrow n+1$) occur at rate μ_1 and loss events ($n \rightarrow n-1$) occur at rate μ_2 . The zeroes indicate that there is no immediate evolutionary path between states i and j : some intermediate state must be used to reach j .

$$Q = \begin{pmatrix} - & r_1 & 0 & 0 \\ r_2 & - & r_1 & 0 \\ 0 & r_2 & - & r_1 \\ 0 & 0 & r_2 & - \end{pmatrix}$$

Create a simplex of unscaled rate parameters for gain ($r[1]$) and loss ($r[2]$) events. The simplex is assigned a flat Dirichlet prior with an initial gain-to-loss rate ratio of 3:1

```
r ~ dnDirichlet( [1,1] )
r.setValue( simplex(3,1) )
```

Create a tridiagonal matrix of transition rates, meaning state i may only transition to states $i-1$ and $i+1$

```
rates := [ [ 0.0, r[1], 0.0, 0.0],
           [ r[2], 0.0, r[1], 0.0],
           [ 0.0, r[2], 0.0, r[1]],
           [ 0.0, 0.0, r[2], 0.0] ]
```

Create the rate matrix

```
Q := fnFreeK(rates)
[ [ -1.5385, 1.5385, 0.0000, 0.0000 ] ,
  [ 0.5128, -2.0513, 1.5385, 0.0000 ] ,
```

```
0.0000, 0.5128, -2.0513, 1.5385 ] ,
0.0000, 0.0000, 0.5128, -0.5128 ] ]
```

Compute the transition probability matrix for a branch length of 0.1.

```
P <- Q.getTransitionProbabilities(rate=0.1)
[ [ 0.861, 0.129, 0.010, 0.001],
  [ 0.043, 0.821, 0.126, 0.010],
  [ 0.001, 0.042, 0.821, 0.136],
  [ 0.000, 0.001, 0.045, 0.954]]
```

Note that $P[1][2] > P[1][3] > P[1][4]$ primarily because those transitions require a minimum of one, two, and three events each. In addition, note that $P[1][2] > P[2][1]$ because $\text{rates}[1][2] > \text{rates}[2][1]$.

7.3 Correlated binary characters

Two characters do not necessarily evolve independently of one another. Take two characters: in plants, the presence or absence of toothed leaf margins (character X) is thought to be ecologically correlated with the presence or absence of leaf lobing (character Y). For a single binary character, there are two states (0 and 1), but there are four for a pair of non-independent binary characters (00, 10, 01, and 11). ? introduced a general framework for modeling the evolution of joint sets of characters. These models require that only one evolutionary event can occur in a moment of time, which is enforced with the 0 terms. In addition, the transition rate that, say, character X goes from 0 to 1 depends on the current value of character Y.

7.3.1 Eight free parameters

In the first case, all possible transitions might be assigned their own parameter.

$$Q = \begin{pmatrix} - & \mu_{00 \rightarrow 10} & \mu_{00 \rightarrow 01} & 0 \\ \mu_{10 \rightarrow 00} & - & 0 & \mu_{10 \rightarrow 11} \\ \mu_{01 \rightarrow 00} & 0 & - & \mu_{01 \rightarrow 11} \\ 0 & \mu_{11 \rightarrow 10} & \mu_{11 \rightarrow 01} & - \end{pmatrix}$$

```
r ~ dnDirichlet( [1,1,1,1,1,1,1,1] )
r.setValue( simplex(1,1,3,3,3,3,1,1) )
```

Create an array of zeroes for the four states (00, 10, 01, 11)

```
for (i in 1:4) {
  for (j in 1:4) {
```

```

      rates[i][j] <- 0.0
    }
  }
}
```

Populate the elements of rates

```

rates[1][2] := r[1] # 00->10
rates[1][3] := r[2] # 00->01
rates[2][1] := r[3] # 10->00
rates[2][4] := r[4] # 10->11
rates[3][1] := r[5] # 01->00
rates[3][4] := r[6] # 01->11
rates[4][2] := r[7] # 11->10
rates[4][3] := r[8] # 11->01
```

Create the rate matrix

```

Q := fnFreeK(rates)
[ [ -0.6667, 0.3333, 0.3333, 0.0000 ] ,
  [ 1.0000, -2.0000, 0.0000, 1.0000 ] ,
  [ 1.0000, 0.0000, -2.0000, 1.0000 ] ,
  [ 0.0000, 0.3333, 0.3333, -0.6667 ] ]
```

Compute the transition probability matrix for a branch length of 0.1.

```

P <- Q.getTransitionProbabilities(rate=0.1)
[ [ 0.938, 0.029, 0.029, 0.003],
  [ 0.088, 0.822, 0.003, 0.088],
  [ 0.088, 0.003, 0.822, 0.088],
  [ 0.003, 0.029, 0.029, 0.938] ]
```

Note that the probability of remaining in state 10 or state 01 is less than the probability of remaining in state 00 or state 11. In this toy example, these probabilities reflect that states 10 and 01 are less evolutionarily stable than states 00 and 11.

7.3.2 Five free parameters

Alternatively, characters X and Y might share state frequencies, π_j , and transition rates $\mu_{ij}^{(k)}$, where i is the starting state for the character undergoing change, j is the ending state, and k is the state of the other character. This results in two stationary frequencies (one free parameter), four transition rates for $0 \rightarrow 1$ and $1 \rightarrow 0$ given that the other character

is in state 0 or state 1 (three free parameters), plus one free parameter to scale the rate matrix.

$$Q = \begin{pmatrix} - & \mu_{01}^{(0)} \pi_1 \pi_0 & \mu_{01}^{(0)} \pi_0 \pi_1 & 0 \\ \mu_{10}^{(0)} \pi_0 \pi_0 & - & 0 & \mu_{01}^{(1)} \pi_1 \pi_1 \\ \mu_{10}^{(0)} \pi_0 \pi_0 & 0 & - & \mu_{01}^{(1)} \pi_1 \pi_1 \\ 0 & \mu_{10}^{(1)} \pi_1 \pi_0 & \mu_{10}^{(1)} \pi_0 \pi_1 & - \end{pmatrix}$$

Assign the stationary frequencies of being in state 0 or state 1 shared by both characters X and Y.

```
pi ~ dnDirichlet([1,1])
pi.setValue( simplex(1,3) )
```

Assign the relative transition rates for gain and loss provided that the other character is in state 0 or 1.

```
r ~ dnDirichlet( [1,1,1,1] )
r.setValue( simplex(1,3,3,1) )
```

Create an array of zeroes for the four states (00, 10, 01, 11)

```
for (i in 1:4) {
  for (j in 1:4) {
    rates[i][j] <- 0.0
  }
}
```

Populate the elements of rates

```
rates[1][2] := r[1] * pi[2] * pi[1] # 00->10
rates[1][3] := r[1] * pi[1] * pi[2] # 00->01
rates[2][1] := r[3] * pi[1] * pi[1] # 10->00
rates[2][4] := r[2] * pi[2] * pi[2] # 10->11
rates[3][1] := r[2] * pi[1] * pi[1] # 01->00
rates[3][4] := r[3] * pi[2] * pi[2] # 01->11
rates[4][2] := r[4] * pi[2] * pi[1] # 11->10
rates[4][3] := r[4] * pi[1] * pi[2] # 11->01
```

Create the rate matrix

```
Q := fnFreeK(rates)
Q
[ [ -0.6333, 0.3167, 0.3167, 0.0000 ] ,
  [ 0.4750, -2.3750, 0.0000, 1.9000 ] ,
  [ 0.4750, 0.0000, -2.3750, 1.9000 ] ,
  [ 0.0000, 0.3167, 0.3167, -0.6333 ] ]
```

Compute the transition probability matrix for a branch length of 0.1.

```
P <- Q.getTransitionProbabilities(0.1)
P
[ [ 0.940, 0.027, 0.027, 0.005],
  [ 0.041, 0.792, 0.003, 0.164],
  [ 0.041, 0.003, 0.792, 0.164],
  [ 0.001, 0.027, 0.027, 0.944] ]
```

Note that this model has a tendency towards state 11.

```
P_10 <- Q.getTransitionProbabilities(10.0)
P_10
[ [ 0.159, 0.105, 0.105, 0.630],
  [ 0.158, 0.105, 0.105, 0.632],
  [ 0.158, 0.105, 0.105, 0.632],
  [ 0.158, 0.105, 0.105, 0.632] ]
```

7.4 Covarion

Covarion models (Tuffley and Steel 1998) capture the possibility that a "hidden" (unobserved or unmeasurable) state causes evolutionary processes to . For example, phylogenetically local clusters of plant lineages transition between herbaceous and woody habits at relatively high rates, and one might want to quantify where these bursts occur (?). While similar to the correlated character model of ?, covarion models do not observe the hidden state that induce the mode-shifts. Instead, covarion models expand the character's state space by a factor of K , and observe the character once for each of the K categories. For example, take a binary character modeled with $K=2$ hidden state classes. The model would treat a character that is observed as being in state 0 as possibly being in either of the $K=2$ classes (0,1) and (0,2). In practice, this is done by setting the likelihood of observing those $0k$ states to equal 1.

The expanded structure of a simple covarion rate matrix with $K=2$ is

$$Q = \left(\begin{array}{cc|cc} - & r_1 q_{01}^{(1)} & s_{12} & 0 \\ r_1 q_{10}^{(1)} & - & 0 & s_{12} \\ \hline s_{21} & 0 & - & r_2 q_{01}^{(2)} \\ 0 & s_{21} & r_2 q_{10}^{(2)} & - \end{array} \right)$$

This form can be reduced to a simpler block-matrix representation

$$Q = \left(\begin{array}{c|c} r_1 Q^{(1)} & s_{12} I \\ \hline s_{21} I & r_2 Q^{(2)} \end{array} \right)$$

where $Q^{(i)}$ is the rate matrix for the i th class, $r_i \in r$ is the clock rate for the i th class, and S is the rate matrix to switch between classes.

```
sr ~ dnDirichlet([1,1])
sr.setValue( simplex(1,2) )
switch_rates := [ [ 0.0, sr[1] ],
                  [ sr[2], 0.0 ] ]
Q_switch := fnFreeK(switch_rates)
```

Create an array of zeroes for the four states (00, 10, 01, 11)

```
cr[1] ~ dnExp(1)
cr[2] ~ dnExp(1)
cr[1].setValue(3)
cr[2].setValue(1)
```

Populate the elements of rates

```
Q_class[1] := fnJC(2)

bf ~ dnDirichlet( [1,1] )
bf.setValue( simplex(1,3) )
Q_class[2] := fnF81( bf )
```

Create the rate matrix

```
Q := fnCovarionRateMatrix(Q=Q_class, switch_rates=Q_switch, clock_rates=cr)

Q
[ [ -1.1126, 0.8901, 0.2225, 0.0000 ] ,
  [ 0.8901, -1.1126, 0.0000, 0.2225 ] ,
  [ 0.4451, 0.0000, -1.0385, 0.5934 ] ,
  [ 0.0000, 0.4451, 0.1978, -0.6429 ] ]
```

Compute the transition probability matrix for a branch length of 0.1.

```
P <- Q.getTransitionProbabilities(0.1)
P <- Q.getTransitionProbabilities(1)
P
[ [ 0.899, 0.080, 0.020, 0.002],
  [ 0.080, 0.899, 0.001, 0.020],
  [ 0.040, 0.003, 0.902, 0.055],
  [ 0.002, 0.041, 0.018, 0.939] ]
```

The rows and columns correspond to (in order): state 0 evolving by $r_1 Q^{(1)}$, state 1 evolving by $r_1 Q^{(1)}$, state 1 evolving by $r_2 Q^{(2)}$, and state 2 evolving by $r_2 Q^{(2)}$. Note that $P[1][2] > P[3][4]$, which is largely due to the fact that $r_1 > r_2$.

References

- Abella, J., P. Montoya, and J. Morales. 2011. Una nueva especie de *Agriarctos* (Ailuropodinae, Ursidae, Carnivora) en la localidad de Nombrevilla 2 (Zaragoza, España). *Estudios Geológicos* 67:187-191.
- Allman, E. S. and J. A. Rhodes. 2008. Identifying evolutionary trees and substitution parameters for the general Markov model with invariable sites. *Mathematical Biosciences* 211:18-33.
- Freudenstein, J. V. 2005. Characters, states and homology. *Systematic Biology* 54:965.
- Huelsenbeck, J. P. and F. Ronquist. 2001. MRBAYES: Bayesian inference of phylogenetic trees. *Bioinformatics* 17:754-755.
- Jukes, T. H. and C. R. Cantor. 1969. Evolution of protein molecules. Pages 21-123 *in* *Mammalian Protein Metabolism* (H. N. Munro, ed.) Academic Press.
- Lartillot, N. and H. Philippe. 2004. A Bayesian mixture model for across-site heterogeneities in the amino-acid replacement process. *Molecular Biology and Evolution* 21:1095-1109.
- Lewis, P. O. 2001. A likelihood approach to estimating phylogeny from discrete morphological character data. *Systematic Biology* 50:913-925.
- Phillips, A. J. 2006. Homology assessment and molecular sequence alignment. *Journal of Biomedical Informatics* 39:18 - 33 *phylogenetic Inferencing: Beyond Biology*.
- Tuffley, C. and M. Steel. 1998. Modeling the covarion hypothesis of nucleotide substitution. *Mathematical Biosciences* 147:63-91.
- Wright, A. M. and D. M. Hillis. 2014. Bayesian analysis using a simple likelihood model outperforms parsimony for estimation of phylogeny from discrete morphological data. *PLoS One* 9:e109210.
- Wright, A. M., G. T. Lloyd, and D. M. Hillis. 2016. Modeling character change heterogeneity in phylogenetic analyses of morphology through the use of priors. *Systematic biology* 65:602-611.

Version dated: March 10, 2017