

# Phylogenetic Inference using RevBayes

*Historical biogeography*

Michael Landis

## 1 Introduction

How did species come to live where they're found today? To answer this, we can leverage phylogenetic and geological information to model species distributions as the outcome of biogeographic processes. These natural processes require some additional considerations, such as how ranges are inherited following speciation events, how geological events might influence dispersal rates, and what factors affect rates of dispersal and extirpation. The major challenge of modeling range evolution is how to translate these natural processes into stochastic processes that remain tractable for inference. This tutorial provides a brief background in some of these models, then describes how to perform Bayesian inference of historical biogeography using RevBayes.

## 2 Dispersal-Extinction-Cladogenesis model

### 2.1 Range characters

Discrete biogeographical models typically rely on presence-absence data, where a species is observed or not observed across multiple discrete areas. For example, say there are three areas: A, B, and C. Say a species is present in areas A and C, then its range equals AC, which can also be encoded into the length-3 bit vector, 101. Bit vectors may also be transformed into (decimal) integers, *e.g.*, the binary number 101 equals the decimal number 5.

$$(\emptyset, A, B, AB, C, AC, BC, ABC) \Leftrightarrow (000, 100, 010, 110, 101, 011, 111) \Leftrightarrow (0, 1, 2, 3, 4, 5, 6, 7)$$

Decimal representation is rarely used in discussion, but it is useful to keep in mind when considering the total number of possible ranges for a species.

### 2.2 Modeling anagenic range evolution

How might we model the dynamics of species range evolution? In this section, we'll cover the Dispersal-Extinction-Cladogenesis model proposed by ?. To begin, we'll focus on anagenesis: evolution that occurs between speciation events within lineages. Since we have discrete characters we'll use the continuous-time Markov chain, which allows us to compute transition probability of a character changing from  $i$  to  $j$  in time  $t$  through matrix exponentiation

$$\mathbf{P}_{i,j}(t) = [\exp \{\mathbf{Q}t\}]_{i,j},$$

where  $\mathbf{Q}$  is the instantaneous rate matrix defining the rates of change between all pairs of characters, and  $\mathbf{P}$  is the transition probability rate matrix. Remember,  $i$  and  $j$  represent different ranges, each of which is encoded as the set of areas occupied by the species. Exponentiation of the rate matrix is powerful because

it integrates over all possible scenarios of character transitions that could occur during  $t$  so long as the chain begins in state  $i$  and ends in state  $j$ .

We can then encode  $\mathbf{Q}$  to reflect the allowable classes of range evolution events with biologically meaningful parameters. We'll take a simple model of range expansion (e.g.  $BC \rightarrow ABC$ ) and range contraction (e.g.  $BC \rightarrow C$ ). (Range expansion may also be referred to as dispersal or area gain and range contraction as extirpation, (local) extinction, or area loss.) The rates in the transition matrix for three areas might appear as

		$\emptyset$	$A$	$B$	$AB$	$C$	$AC$	$BC$	$ABC$
$\mathbf{Q} =$	$\emptyset$	—	0	0	0	0	0	0	0
	$A$	$e_A$	—	0	$d_{AB}$	0	$d_{AC}$	0	0
	$B$	$e_B$	0	—	$d_{BA}$	0	0	$d_{BC}$	0
	$AB$	0	$e_A$	$e_B$	—	0	0	0	$d_{AC} + d_{BC}$
	$C$	$e_C$	0	0	0	—	$d_{CA}$	$d_{CB}$	0
	$AC$	0	$e_C$	0	0	$e_A$	—	0	$d_{AB} + d_{CB}$
	$BC$	0	0	$e_C$	0	$e_B$	0	—	$d_{BA} + d_{CA}$
	$ABC$	0	0	0	$e_C$	0	$e_B$	$e_A$	—

where  $e = (e_A, e_B, e_C)$  are the (local) extinction rates per area, and  $d = (d_{AB}, d_{AC}, d_{BC}, d_{CB}, d_{CA}, d_{BA})$  are the dispersal rates between areas. Notice that the sum of rates leaving state  $\emptyset$  is zero, meaning any species that loses all areas in its range remains permanently extinct.

?

For the three-area DEC rate matrix above, what is the rate of leaving state AC in terms of dispersal and extinction parameters?

Note the rate of more than one event occurring simultaneously is zero, so a range must expand twice by one area in order to expand by two areas.

?

What series of transition events might explain a lineage evolving from range  $ABC$  to range  $A$ ? From range  $AB$  to range  $C$ ?

Of course, this model can be specified for more than three areas.

?

Imagine a DEC rate matrix with four areas,  $ABCD$ . What would be the dispersal rate for  $Q_{BC,BCD}$ ? How many states does a DEC rate matrix with four areas have? What is the relationship between the number of areas and the number of states under the DEC model?

Let's consider what happens to the size of  $\mathbf{Q}$  when the number of areas,  $N$ , becomes large. For three areas,  $\mathbf{Q}$  is size  $8 \times 8$ . For ten areas,  $\mathbf{Q}$  is size  $2^{10} \times 2^{10} = 1024 \times 1024$ , which approaches the largest size matrices that can be exponentiated in a practical amount of time. For twenty areas,  $\mathbf{Q}$  is size  $2^{20} \times 2^{20} \approx 10^6 \times 10^6$  and exponentiation is not viable.

## 2.3 Modeling cladogenic range evolution

Cladogenesis describes evolutionary change accompanying speciation. Daughter species are not expected to inherit their ancestral range identically in general. For each internal node in the reconstructed tree, one of two cladogenic events can occur: sympatry or allopatry. Say the range of a species is  $A$  the moment before speciation occurs at an internal phylogenetic node. Since the species range is size one, both daughter lineages necessarily inherit the ancestral species range ( $A$ ). In DEC parlance, this is called a *narrow*

*sympatry* event.

Now suppose the ancestral range is  $ABC$ . Under *subset sympatric cladogenesis*, one lineage identically inherits the ancestral species range,  $ABC$ , while the other lineage inherits only a single area, i.e. only  $A$  or  $B$  or  $C$ . For *widespread sympatric cladogenesis*, both lineages inherit the ancestral range,  $ABC$ . Under *allopatric cladogenesis*, the ancestral range is split evenly among daughter lineages, e.g. one lineage may inherit  $AB$  and the other inherits  $C$ .

For an excellent overview of described state transitions for cladogenic events, see ?.

**[?] Given the state is  $AB$  before cladogenesis, and allowing subset sympatry, widespread sympatry, and allopatry, what are the 7 possible states in the daughter lineages after cladogenesis?**

The probabilities of anagenic change along lineages must account for all combinations of starting states and ending states. For 3 areas, there are 8 states, and thus  $8 \times 8 = 64$  probability terms for pairs of states. For cladogenic change, we need transition probabilities for all combinations of states before cladogenesis, after cladogenesis for the left lineage, and after cladogenesis for the right lineage. Like above, for three areas, there are 8 states, and  $8 \times 8 \times 8 = 512$  cladogenic probability terms.

**[?] For three areas, there are three narrow, four widespread, 18 subset sympatric events and 12 allopatric cladogenesis events. What proportion of terms in the cladogenesis matrix are zero?**

The DEC model ignores speciation events hidden by extinction or incomplete taxon sampling. The probability of cladogenesis and local extinction events would ideally be linked to a birth-death process, as it is in the GeoSSE model (?). Unfortunately, since the numerical method for SSE models scale poorly, and DEC models remain the only option when the geography has more than two or three areas. For more than ten areas, data augmentation may be used to infer ancestral ranges, as described in Section ??.

The rest of this section will describe how to run a simple DEC analysis using RevBayes.

### 2.3.1 Specifying a simple DEC model

We'll use the primate dataset with 23 taxa. To keep the model simple, we'll discretize their ranges into just three areas: the New World (A), Africa (B), and Eurasia (C). For simplicity, we'll assume their phylogeny is time-calibrated, errorless, and fixed.

Create some `String` variables for file handling,

```
data_fn = "data/primates_bg_n3.tsv"
tree_fn = "data/primates.tree"
out_fn = "output/bg_same"
```

then read in our character data, using `type="Bitset"` to indicate ranges are encoded with bits, e.g. 011

```
data = readCharacterDataDelimited(file=data_fn, type="Bitset")
```

and our tree

```
psi <- readTrees(tree_fn)[1]
```

Next, compute the number of states from the number of areas

```
n_areas = 3
n_states = floor(2^n_areas)
```

Declare index variables for our move vectors for future use

```
mi = 0
```

Now, we'll begin to construct the rate matrix for anagenic events. First create a matrix, 8-by-8 in size, initialized with all zeroes

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

Now we need to populate the non-zero rate matrix elements, which are in terms of dispersal and extinction rates. We'll use one dispersal rate and one extinction rate for this tutorial, and explore more complex models in later sections. For later reference, this will be called the “same rate” model.

First, create a extinction rate parameter and assign it a scale move

```
r_e ~ dnExponential(10.0)
mv[++mvi] = mvScale(r_e, weight=5)
```

Before assigning the rates to the rate matrix, we'll create a vector to hold the per-area extinction rates

```
for (i in 1:n_areas) {
    e[i] := r_e
}
```

Now create the dispersal rate and scale move

```
r_d ~ dnExponential(10.0)
mv[++mvi] = mvScale(r_d, weight=5)
```

then assign the between-area dispersal rates as determined by  $r_d$

```
for (i in 1:n_areas) {
    for (j in 1:n_areas) {
        d[i][j] <- 0.
        if (i != j) {
            d[i][j] := r_d
        }
    }
}
```

Although the DEC rate matrix can be easily constructed by typing

```
q := fnDECRateMatrix(d,e)
```

manually encoding the structure of the DEC rate matrix illuminates the relationship between dispersal rates, extinction rates, area gain rates, and area loss rates. To start, we'll populate the non-zero rate matrix elements. Rates are indexed by the natural number value of the range (plus one), e.g. the range spanning Eurasia and Africa is coded as 011, which is state 4.

First assign the extinction (range loss) rates

```
r[2][1] := e[1]          # 100 -> 000 : Extirpate in area 1
r[3][1] := e[2]          # 010 -> 000 : Extirpate in area 2
r[4][2] := e[2]          # 011 -> 001 : Extirpate in area 2
r[4][3] := e[3]          # 011 -> 010 : Extirpate in area 3
r[5][1] := e[3]          # 001 -> 000 : Extirpate in area 3
r[6][2] := e[1]          # 101 -> 001 : Extirpate in area 1
r[6][5] := e[3]          # 101 -> 100 : Extirpate in area 3
r[7][3] := e[1]          # 110 -> 010 : Extirpate in area 1
r[7][5] := e[2]          # 110 -> 100 : Extirpate in area 2
r[8][4] := e[1]          # 111 -> 011 : Extirpate in area 1
r[8][6] := e[2]          # 111 -> 101 : Extirpate in area 2
r[8][7] := e[3]          # 111 -> 110 : Extirpate in area 3
```

then the dispersal (range gain) rates

```
r[2][4] := d[3][2]      # 001 -> 011 : Disperse from area 3 to 2
r[2][6] := d[3][1]      # 001 -> 101 : Disperse from area 3 to 1
r[3][4] := d[2][3]      # 010 -> 011 : Disperse from area 2 to 3
r[3][7] := d[2][1]      # 010 -> 110 : Disperse from area 2 to 1
r[5][6] := d[1][3]      # 100 -> 101 : Disperse from area 1 to 3
r[5][7] := d[1][2]      # 100 -> 110 : Disperse from area 1 to 2
r[4][8] := d[2][1] + d[3][1] # 011 -> 111 : Disperse from area 2 to 1 and from 3 to 1
r[6][8] := d[1][2] + d[3][2] # 101 -> 111 : Disperse from area 1 to 2 and from 3 to 2
r[7][8] := d[1][3] + d[2][3] # 110 -> 111 : Disperse from area 1 to 3 and from 2 to 3
```

Show the value of `r` and compare it to the matrix in Section 2.2.

Of course we did not need to declare `d` and `e` to assign `r`, but we'll see these intermediate variables act as a template expose the structure of `r` for modification.

So far, we only have the desired parameterization of the rate matrix, but we still haven't created a rate matrix function. Converting the vector-of-vectors, `r`, into a simplex allows us to use existing rate matrix functions.

First, we'll convert `r` into a one-dimensional vector, skipping the diagonal elements.

```
k = 1
for (i in 1:n_states) {
    for (j in 1:n_states) {
        if (i != j) {
            er_nat[k++] := r[i][j]
        }
    }
}
```

Finally, normalize `er_nat` using a simplex, then pass the resulting exchangeability rates as arguments into the rate matrix function, `q`.

```
er := simplex(er_nat)
q := fnFreeK(er)
```

This yields the desired three-area DEC rate matrix modeling anagenic character change.

In contrast, cladogenic event probabilities are given by a transition probability matrix and do not require a rate matrix. First, we will create a vector of prior weights on cladogenesis events. Here, we assign a flat prior to all cladogenic events

```
widespread_sympatry_wt <- 1.0
subset_sympatry_wt   <- 1.0
allopatry_wt          <- 1.0
clado_prior           <- [ widespread_sympatry_wt, subset_sympatry_wt, allopatry_wt ]
```

then create the distribution over cladogenic event types and add its MCMC move

```
clado_type ~ dnDirichlet(clado_prior)
mv[++mvi] = mvSimplexElementScale(clado_type, alpha=10, weight=5)
```

To give the simplex elements descriptive names when monitored, assign the values to deterministic nodes

```
widespread_sympatry := clado_type[1]
subset_sympatry      := clado_type[2]
allopatry             := clado_type[3]
```

Then create the cladogenic transition probability matrix, which assigns probabilities to cladogenic event classes according to `clado_type_prob`

```
clado_prob := fnCladoProbs(clado_type, n_areas, 2)
```

Add a parameter for a biogeographical clock, which scales the overall rate of range evolution. As a prior, an exponential distribution with rate 10 generates one dispersal or extinction event per 10 million years.

```
clock_bg ~ dnExponential(10)
mv[++mvi] = mvScale(clock_bg, weight=5)
```

Finally, all our model components are encapsulated in the `dnPhyloCTMCClado` distribution, which is similar to `dnPhyloCTMC` except specialized to integrate over cladogenic events. Although this dataset has three areas, it is recognized single character with states valued from 1 to  $2^3$ , hence `nSites=1`.

```
m ~ dnPhyloCTMCClado( tree=psi, Q=q, cladoProbs=clado_prob, branchRates=clock_bg, nSites =1, type="NaturalNumbers" )
```

The remaining tasks should be familiar by now, so we can proceed briskly. Attach the observed ranges to the model.

```
m.clamp(data)
```

Compose the model.

```
mdl = model(m)
```

Add the monitors. (The `mnJointConditionalAncestralState` monitor will be described in the next section.)

```
mn[1] = mnScreen(clock_bg, d[1][2], d[1][3], d[2][1], d[2][3], d[3][1], d[3][2], e[1], e [2], e[3], widespread_sympatry, subset_sympatry, allopatry, printgen=1000)
mn[2] = mnFile(clock_bg, d[1][2], d[1][3], d[2][1], d[2][3], d[3][1], d[3][2], e[1], e [2], e[3], widespread_sympatry, subset_sympatry, allopatry, file=out_fn+".params.txt ")
mn[3] = mnJointConditionalAncestralState(tree=psi, ctmc=m, filename=out_fn+".states.txt ", type="NaturalNumbers", printgen=10, withTips=true, withStartStates=true)
```

Create the MCMC object, and run the chain after burn-in.

```
ch = mcmc(mv,mn,mdl)
ch.burnin(1000, 10)
ch.run(10000)
```

### 2.3.2 Per-area rates

Biologically, local extinction events probably do not occur at equal rates across all areas, as done above. Ecological factors, geographical distances, etc. might cause these parameters to be weakly correlated or completely uncorrelated. Dispersal rates, also, might not be the same between pairs of areas, or even symmetric depending on the direction of dispersal. Rather than constraining all events of a type to share a common rate, instead you might give each area its own extinction parameter

```
for (i in 1:3) {
    e[i] ~ dnExponential(10.0)
    mv[++mvi] = mvScale(e[i], weight=5)
}
```

or give each ordered pair of areas it's own dispersal rate

```
for (i in 1:3) {
    for (j in 1:3) {
        d[i][j] <- 0.0
        if (i != j) {
            d[i][j] ~ dnExponential(10.0)
            mv[++mvi] = mvScale(d[i][j], weight=5)
        }
    }
}
```

Note that you don't need the global dispersal rate `r_d` and extinction rate `r_e` anymore and you can remove the variable from your analysis. RevBayes might give you an error message if you have left them in.

## 2.4 Exercises

Exercises are independent of each other, except for Exercises 3a and 3b.

- 1) Widespread sympatric speciation is thought to be evolutionarily rare. Set the Dirichlet prior on cladogenic event types to heavily disfavor these events. Using Tracer, describe how changing the cladogenesis prior affects the extinction rate when compared with the “common rate” model.
- 2) Modifying the RevBayes script, parameterize the rate matrix so ranges may not grow beyond two areas in size. You may use `q_test := fnDECRateMatrix(e,d,2)` to confirm your results, which will give the same rate structure (and rates, up to a rescaling constant).
- 3a) Saving your commands to a file, create a script to produce the “per-area rate” model.
- 3b) Determine if the data support the “common rate” model over the “per-area rate” model. Use the stepping stone method to compute marginal likelihoods, which will let you compute Bayes factors for model selection.

## 2.5 (Advanced) Joint inference of phylogeny and historical biogeography

This is a slightly more advanced question that you may want to skip if you will always use known, fixed trees in your analyses. Using what you learned in this tutorial and the CTMC tutorial, perform a joint analysis of molecular and biogeographic evolution for primates.

Start by loading the sequence data matrix specified in `data/primates_cytb.nex`.

```
seqData <- readDiscreteCharacterData("data/primates_cytb.nex")
```

We need to get some useful variables from the data so that we will be able to specify the tree prior below. These variables are the number of tips, the number of nodes and the names of the species which we all can query from the continuous character data object.

```
numTips = seqData.ntaxa()
names = seqData.names()
numNodes = numTips * 2 - 1
```

Instead of having a fixed tree as in the previous, we should now define a *random* tree. We use a birth death prior with prior distributions on the **speciation** rate and **extinction** rate.

```
speciation ~ dnExponential(10.)
extinction ~ dnExponential(10.)
moves[++mvi] = mvScale(speciation, lambda=1, tune=true, weight=3.0)
moves[++mvi] = mvScale(extinction, lambda=1, tune=true, weight=3.0)
```

The phylogeny that we used are obviously not a complete sample of all the species and you should take the incomplete sampling into account. We will simply use an empirical estimate of the fraction of species which we included in this study. For more information about incomplete taxon sampling see `?` and `?`.

```
sampling_fraction <- 23 / 270 # 23 out of the ~ 270 primate species
```

Now we are able to specify of tree variable **psi** which is drawn from a constant rate birth-death process. We will condition the age of the tree to be 75 million years old which is approximately the crown age of primates, although this estimate is still debated. We only condition here on the crown age for simplicity because we do not use any other fossil calibration.

```
psi ~ dnBDP(lambda=speciation, mu=extinction, rho=sampling_fraction, rootAge=75, nTaxa=numTips,
            names=names)
```

Note that, here, we do not have included any fossil information: we are merely doing *relative* dating.

The first moves on the tree which we specify are moves that change the node ages. The first move randomly picks a subtree and rescales it, and the second move randomly pick a node and uniformly proposes a new node age between its parent age and oldest child's age.

```
moves[++mvi] = mvSubtreeScale(psi, weight=5.0)
moves[++mvi] = mvNodeTimeSlideUniform(psi, weight=10.0)
```

We also need moves on the tree topology to estimate the phylogeny. The two moves which you use are the nearest-neighbor interchange (NNI) and the fixed-nodeheight-prune-and-regraft (FNPR) (?).

```
moves[++mvi] = mvNNI(psi, weight=5.0)
moves[++mvi] = mvFNPR(psi, weight=5.0)
```

In the next step we set up the substitution model. First, we create a substitution model, just like what you probably did in previous tutorial (*e.g.*, RB\_CTMC\_Tutorial). In a first step, we will use a GTR+Gamma model. We can use a flat Dirichlet prior density on the exchangeability rates **er\_mol** and the base frequencies **pi\_mol**.

```
er_mol_prior <- v(1,1,1,1,1,1)
er_mol ~ dnDirichlet(er_mol_prior)
pi_mol_prior <- v(1,1,1,1)
pi_mol ~ dnDirichlet(pi_mol_prior)
```

Now add the simplex scale move one each the exchangeability rates **er\_mol** and the stationary frequencies **pi\_mol** to the moves vector:

```
moves[++mvi] = mvSimplexElementScale(er_mol, weight=15)
moves[++mvi] = mvSimplexElementScale(pi_mol, weight=5)
```

We can finish setting up this part of the model by creating a deterministic node for the GTR instantaneous-rate matrix **Q\_mol**. The **fnGTR()** function takes a set of exchangeability rates and a set of base frequencies to compute the instantaneous-rate matrix used when calculating the likelihood of our model.

```
Q_mol := fnGTR(er_mol, pi_mol)
```

The next part of the substitution process is the rate variation among sites. We will model this using the commonly applied 4 discrete gamma categories which only have a single parameter **alpha**. Let us specify the rate of **alpha** to 0.05 (thus the mean will be 20.0).

```
alpha_prior <- 0.05
```

Then create a stochastic node called **alpha** with an exponential prior:

```
alpha ~ dnExponential(alpha_prior)
```

Initialize the **gamma\_rates** deterministic node vector using the **fnDiscretizeGamma()** function with 4 bins:

```
gamma_rates := fnDiscretizeGamma( alpha, alpha, 4 )
```

The random variable that controls the rate variation is the stochastic node **alpha**. We will apply a simple scale move to this parameter.

```
moves[++mvi] = mvScale(alpha, weight=2.0)
```

This finishes the substitution process part of the model.

Then next part of the model is the clock model. Here we need a clock model because we work on a time tree. We use an exponential distribution with expectation 0.1.

```
clock_mol ~ dnExponential(10)
moves[++mvi] = mvScale(clock_mol, lambda=1, tune=true, weight=2.0)
```

Introduce a common prior for the molecular and biogeographical clocks.

```
clock_scale_bg ~ dnGamma(2.0,2.0)
clock_bg      := clock_mol * clock_scale_bg
```

Remember that you need to call the **PhyloCTMC** constructor to include the new site-rate parameter:

```
seq ~ dnPhyloCTMC(tree=psi, Q=Q_mol, siteRates=gamma_rates, branchRates=clockRate,
type="DNA")
```

Finally we need to attach the molecular sequence data to our model.

```
seq.clamp(seqData)
```

We are essentially done now. We only need to add a new monitor for the tree so that we can monitor and build the maximum a posteriori tree later.

```
monitors[3] = mnFile(filename="output/biogeography_DEC_joint.trees", printgen=100, separator =
TAB, psi)
```

Create the range evolution model as before, being sure to use the same `psi` and `clock_bg` from above.

The remaining challenge is to compose both submodels together using `model` and creating an MCMC object with `mcmc`. Good luck!

## 3 Epoch models and ancestral range reconstruction

### 3.1 Ancestral range reconstruction

From the previous section, we created an ancestral state monitor by

```
mn[3] = mnJointConditionalAncestralState(tree=psi, ctmc=m, filename=out_fn+".states.txt"
", type="NaturalNumbers", printgen=10, withTips=true, withStartStates=true)
```

In this section, we'll use the output from running the DEC analysis assuming all each ordered pair of areas has its own dispersal parameter, and each area has its own extinction parameter. To generate these results, run

```
source("RevBayes_scripts/biogeography_DEC_diff.Rev")
```

The joint-conditional ancestral state monitor samples the joint distribution of ancestral states every `printgen` iterations over the entire tree, conditioning on the observed tip states. Due to cladogenesis, the state before speciation and the states inherited after speciation may differ, which we monitor setting `withStartStates=true`. For convenience we record the tip states using `withTips=true` though these are known through the input data.

The resulting file, `bg_diff.states.txt`, appears as

Iteration	start_0	end_0	start_1	end_1	start_2	end_2	start_3	end_3	...
0	2	2	2	4	2	2	2	6	...
10	2	2	4	4	2	2	2	6	...
20	2	2	2	4	2	2	6	6	...
...									

where columns give the integer-valued range corresponding to either the start or end of a branch leading to a particular node, and each row corresponds to a single sample drawn from the joint distribution of ancestral states. For example, the lineage leading to the node indexed 1 starts in state 2 (Africa) and ends in state 4 (Eurasia) at iteration 20. In Section ??, we will look at how stochastic mappings—the completely realized biogeographic history—may be analysed for data-augmented models. If you wish to write your own scripts to analyse the ancestral range reconstructions, know that each node's index is recorded when writing a tree's Newick string to file so states may be mapped on to the tree.

### 3.2 Data exploration with Phylowood

To interpret the biogeographical history of primates, we will generate a Phylowood (<http://mlandis.github.io/phylowood>) animation. First, we'll create variables for the relevant files

```
state_fn = "output/bg_diff.states.txt"
atlas_fn = "data/earth3.still.atlas.txt"
phw_fn = "output/bg_diff.phw.txt"
```

then create the animation file

```
convertToPhylowood(treefile=tree_fn, geofile=atlas_fn, statefile=state_fn, outfile=
    phw_fn, burnin=0., chartype="NaturalNumbers", bgtype="Range")
```

This file summarizes the MCMC output from a RevBayes biogeographical analysis as a Nexus-formatted file, which is used by Phylowood to generate interactive animations to explore biogeographic reconstructions.

- Open <http://mlandis.github.io/phylowood>.
- Drag and drop `./output/bg_diff.phw.txt` into the text field.

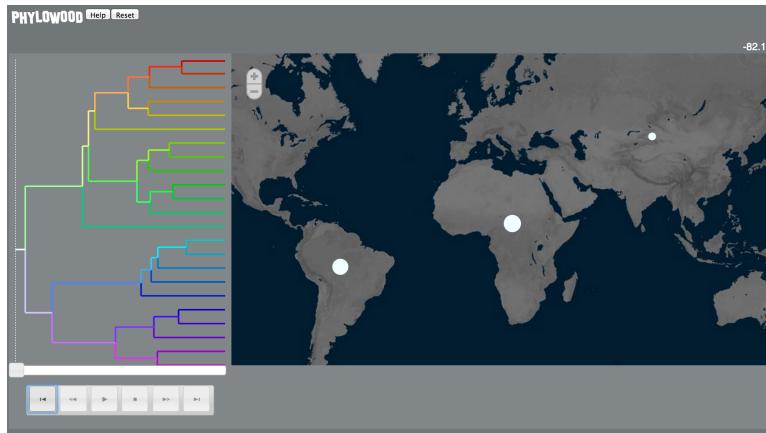


Figure 1: Phylowood frame showing posterior ancestral range of root node.

- Click the Play button to view the animation.

There are three control panels to help you filter data: the media panel, the map panel, and the phylogeny panel. The media buttons correspond to Beginning, Slow/Rewind, Play, Stop, Fast Forward, Ending (from left to right). The animation will play the timeframe corresponding to the slider.

- Drag the slider to the right (the present).

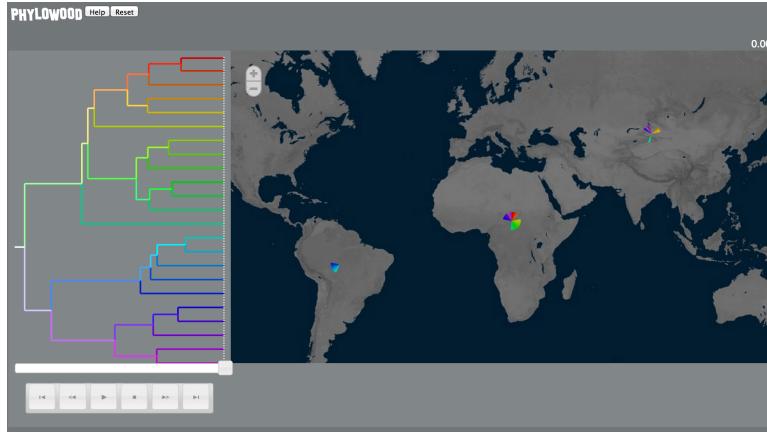


Figure 2: Phylowood frame showing distribution of extant taxon ranges.

- Pan and zoom around the map.

Marker colors correspond to the phylogenetic lineages in the phylogeny panel. Markers are split into slices and (loosely) sorted phylogenetically, so nearby slices are generally closely related. At divergence events, a marker's radius is proportional to the marginal posterior probability the node was present in the area at that time. Between divergence events, marker's radius is simply an interpolation of the values at the two endpoints. Some information about geological constraints and cladogenic events is lost.

- Mouseover an area to learn which lineage it belongs to and its presence probability.

Since it's difficult to see how specific clades evolve with so many taxa, Phylowood offers two ways to filter taxa from the animation. We call the set of a lineage, all its ancestral lineages towards the root, and all descendant lineages a phylogenetic heritage. The root's heritage is the entire clade. A leaf node's heritage is a path from the tip to the root.

- Mouseover a lineage to temporarily highlight the lineage's heritage. Remove the mouseover to remove the highlight effect.

The highlight effect is temporary and quickly allows you to single out lineages of interest during animation. Phylowood also offers a masking effect that persists until an unmask command is issued.

- Double-click the white root branch to mask the root node's heritage (all lineages). Single click a lineage to unmask that lineage's heritage.

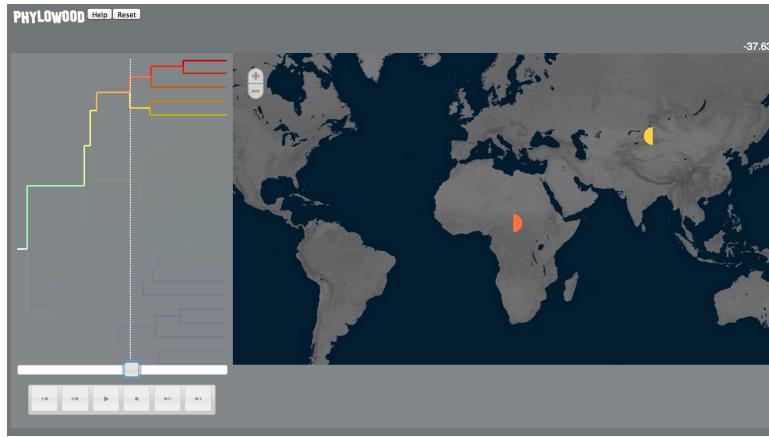


Figure 3: Phylowood frame highlighting the ancestral range for the MRCA extant lorises.

Now that the masking effects are in place, you’re free to interact with other map components. In addition, the area of marker sizes is only distributed among unmasked lineages.

Visit <https://github.com/mlandis/phylowood/wiki> to learn more about Phylowood.

Phylowood is useful to understand whether your model generates sensible reconstructions. Keep in mind the animation and inference both use modern geographies, and that dispersal rates are not modeled to depend on geographical features. Notice the primate MRCA is widespread in two areas, the New World and Africa. The root age of the primate tree is about 80Mya, which is about 40Mya after East and West Gondwanaland split to separate modern Africa and South America, so vicariance did not give rise to New World Monkeys. This is corroborated by examining the fossil record, where the first appearance of primates in South America is about 35Mya. The standing theory is primates disperse from Africa to South America in a rare “sweepstakes.” The naive primate biogeographic model described so far should produce more realistic reconstructions by taking continental drift into account.

### 3.3 Epoch models

To model the effects of geography, we will use *epoch* models. Rather than assuming the evolutionary process is constant with respect to time, it assumes the process is piecewise-constant. For example, Africa (along with all of Gondwanaland) split from Eurasia (and all of Laurasia) around 180 Mya. Africa merged with Eurasia only about 50 Mya. It is reasonable to expect that dispersal rates between Africa and Eurasia were lower before 50 Mya, and higher after 50 Mya until present.

To model this, we will specify one rate matrix that describes anagenic dispersal and extinction processes before for ages  $t \geq 50$  and a second rate matrix for the process operating after  $0 \leq t < 50$ .

To proceed, first we will read in an **Atlas** file, which fully describes the geography per time interval. For more on the **Atlas** format, see Section XXX.

```
# read the atlas
atlas <- readAtlas("data/earth3.drift.atlas.txt")
n_epochs = atlas.nEpochs()
times <- atlas.epochTimes()
```

This atlas contains two epochs, each with three areas, and a single breakpoint at age 50 (50Mya). For example, the area for Africa in the second epoch from 50Mya – present reads

```
{
  "latitude": 10.0,
  "longitude": 20.0,
  "dispersalValues": [ 0.1, 0.0, 1.0 ],
  "extinctionValues": [ 1.0 ],
  "name": "Africa"
}
```

The `dispersalValues` values share the ordering of the areas, and each area reports a value regarding it's relationship to other areas. Here, 0.1 in position 1 indicates the Atlantic Ocean obstructs dispersal from Africa to the Americas, while the 1.0 in position 3 indicates African primates may easily disperse into Eurasia by way of the Arabian Peninsula. (Note 0.0 in position 2 is the dispersal value from Africa to itself, so the value is arbitrary.)

To extract these values in matrix form

```
d_prior <- atlas.getValues("dispersal")
```

where, for example, the dispersal rate more than 50 Myr in the past (epoch 1) from Africa (area 2) to Eurasia (area 3) is accessed by

```
d_prior[1][2][3]
```

To use these empirical priors to rescale our naive priors

```
for (t in 1:n_epochs) {
    for (i in 1:n_areas) {
        for (j in 1:n_areas) {
            if (i != j) {
                d_raw[t][i][j] ~ dnExp( 10. )
                mv[mvi++] = mvScale(d_raw[t][i][j], weight=2)
                d[t][i][j] := d_raw[t][i][j] * abs(d_prior[t][i][j])
            } else {
                d[t][i][j] <- abs(0.)
            }
        }
    }
}
```

We'll make no strong prior assumptions about epochs or areas differentially affecting extinction rates.

```
# set up extinction per epoch
e_prior <- atlas.getValues("extinction")
for (t in 1:n_epochs) {
    for (i in 1:n_areas) {
        e_raw[t][i] ~ dnExp( 10. )
        mv[mvi++] = mvScale(e_raw[t][i], weight=2)
        e[t][i] := e_raw[t][i] * abs(d_prior[t][i][1])
    }
}
```

Now we have dispersal rates and extinction rates per epoch in the matrices `d` and `e`.

Extant primate ranges do not appear capable of spanning all three areas simultaneously, so we might constrain the range evolution process to only allow ranges of size two

```
rangeSize <- simplex(0,1,1,0)
```

where the first and last elements are set to zero, so ranges cannot be size 0 or size 3. By calling `fnDECRoot` with the `rangeSize` parameter, we can also force the MRCA range size to be a single area.

```
rf := fnDECRoot( rep(1,n_states), rangeSize=simplex(0,1,0,0) )
```

Next, we'll create a rate matrix for each epoch

```
for (t in 1:n_epochs) {
    rates[t] ~ dnGamma(2,2)
    mv[mvi++] = mvScale(rates[t], weight=2)
    q[t] := fnDECRateMatrix(d[t], e[t], rangeSize)
}
```

with `rates` being an epochal rate multiplier, each with mean 1. Notice, we construct `q[t]` for each epoch `t` in `n_epochs`, using the epoch's dispersal values `d[t]` and extinction values `e[t]`. Finally, we wrap our vector of rate matrices, `q`, along with the epoch boundary times, `times`, and our epochal rates, `rates`

```
q_epoch := fnEpoch( Q=q, times=times, rates=rates )
```

Parameters like `clado_prob` and `clock_bg` still need to be created, as they were in the previous section.

Otherwise, the model is created as before, except passing `q_epoch` into `dnPhyloCTMCClado`.

```
m ~ dnPhyloCTMCClado( tree=psi, Q=q_epoch, rootFrequencies=rf, cladoProbs=clado_prob,
  branchRates=clock_bg, nSites=1, type="NaturalNumbers" )
m.clamp(data)
mdl = model(m)
```

This biologically and geographically informed model produces a more realistic range reconstruction

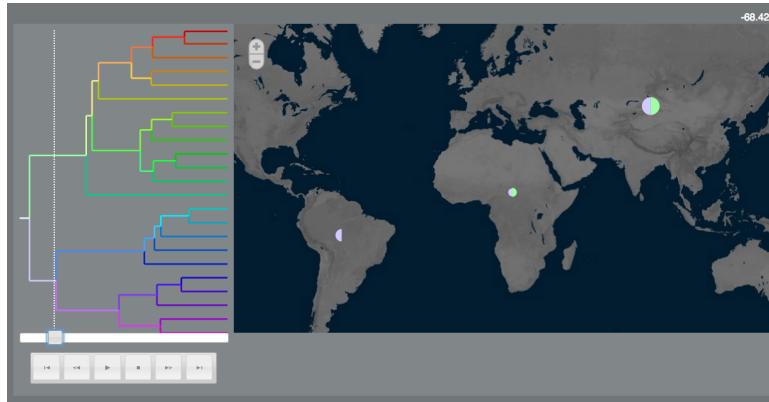


Figure 4: Phylowood frame showing Asian origin of primates, with subsequent dispersal into Africa and Americas.

This example model may be run by typing

```
RevBayes_scripts/biogeography_epoch.Rev}
```

which also produces the above Phylowood animation file.

## References

## References

Version dated: July 10, 2016