

# Phylogenetic Inference using RevBayes

## *Basic Diversification Rate Estimation*

Sebastian Höhna and Tracy Heath

## 1 Overview: Diversification Rate Estimation

Models of speciation and extinction are fundamental to any phylogenetic analysis of macroevolutionary processes. A prior describing the distribution of speciation events over time is critical to estimating phylogenies with branch lengths proportional to time. Moreover, stochastic branching models allow for inference of speciation and extinction rates. These inferences allow us to investigate key questions in evolutionary biology.

Similarly, diversification-rate parameters are also included as nuisance parameters of other phylogenetic models—*i.e.*, where these diversification-rate parameters are not of direct interest. For example, many methods for estimating species divergence times—such as BEAST (Drummond et al. 2012), MrBayes (Ronquist et al. 2012), and RevBayes (Höhna et al. 2015)—implement ‘relaxed-clock models’ that include a constant-rate birth-death branching process as a prior model on the distribution of tree topologies and node ages. Although the parameters of these ‘tree priors’ are not typically of direct interest, they are nevertheless estimated as part of the joint posterior probability distribution of the relaxed-clock model, and so can be estimated simply by querying the corresponding marginal posterior probability densities. In fact, this may provide more robust estimates of the diversification-rate parameters, as they accommodate uncertainty in the other phylogenetic-model parameters (including the tree topology, divergence-time estimates, and the other relaxed-clock model parameters).

### 1.1 Types of Hypotheses for Estimating Diversification Rates

Many evolutionary phenomena entail differential rates of diversification (speciation – extinction); *e.g.*, adaptive radiation, diversity-dependent diversification, key innovations, and mass extinction. The specific study questions regarding lineage diversification may be classified within three fundamental categories of inference problems. Admittedly, this classification scheme is somewhat arbitrary, but it is nevertheless useful, as it allows users to navigate the ever-increasing number of available phylogenetic methods. Below, we describe each of the fundamental questions regarding diversification rates.

**(1) Diversification-rate through time estimation** *What is the (constant) rate of diversification in my study group?* The most basic models estimate parameters of the stochastic-branching process (*i.e.*, rates of speciation and extinction, or composite parameters such as net-diversification and relative-extinction rates) under the assumption that rates have remained constant across lineages and through time; *i.e.*, under a constant-rate birth-death stochastic-branching process model. Extensions to the (basic) constant-rate models include diversification-rate variation through time. First, we might ask whether there is evidence of an episodic, tree-wide increase in diversification rates (associated with a sudden increase in speciation rate and/or decrease in extinction rate), as might occur during an episode of adaptive radiation. A second question asks whether there is evidence of a continuous/gradual decrease in diversification rates

through time (associated with decreasing speciation rates and/or increasing extinction rates), as might occur because of diversity-dependent diversification (*i.e.*, where competitive ecological interactions among the species of a growing tree decrease the opportunities for speciation and/or increase the probability of extinction). A final question in this category asks whether our study tree was impacted by a mass-extinction event (where a large fraction of the standing species diversity is suddenly lost).

**(2) Diversification-rate variation across branches estimation** *Is there evidence that diversification rates have varied significantly across the branches of my study group?* Models have been developed to detect departures from rate constancy across lineages; these tests are analogous to methods that test for departures from a molecular clock—*i.e.*, to assess whether substitution rates vary significantly across lineages. These models are important for assessing whether a given tree violates the assumptions of other inference methods. Furthermore, these models are important to answer questions such as: *What are the branch-specific diversification rates?*; and *Have there been significant diversification-rate shifts along branches in my study group, and if so, how many shifts and along which branches?*

**(3) Character-dependent diversification-rate estimation** *Are diversification rates correlated with some variable in my study group?* Character-dependent diversification-rate models aim to identify overall correlations between diversification rates and organismal features (binary and multi-state discrete morphological traits, continuous morphological traits, geographic range, etc.). For example, one can hypothesize that a binary character, say if an organism is herbivorous/carnivorous or self-compatible/self-incompatible, impact the diversification rates. Then, if the organism is in state 0 (*e.g.*, is herbivorous) it has a lower (or higher) diversification rate than if the organism is in state 1 (*e.g.*, carnivorous).

## 2 Models

We begin this section with a general introduction to the stochastic birth-death branching process that underlies inference of diversification rates in **RevBayes**. This primer will provide some details on the relevant theory of stochastic-branching process models. We appreciate that some readers may want to skip this somewhat technical primer; however, we believe that a better understanding of the relevant theory provides a foundation for performing better inferences. We then discuss a variety of specific birth-death models, but emphasize that these examples represent only a tiny fraction of the possible diversification-rate models that can be specified in **RevBayes**.

### 2.1 The birth-death branching process

Our approach is based on the *reconstructed evolutionary process* described by [Nee et al. \(1994\)](#); a birth-death process in which only sampled, extant lineages are observed. Let  $N(t)$  denote the number of species at time  $t$ . Assume the process starts at time  $t_1$  (the ‘crown’ age of the most recent common ancestor of the study group,  $t_{\text{MRCA}}$ ) when there are two species. Thus, the process is initiated with two species,  $N(t_1) = 2$ . We condition the process on sampling at least one descendant from each of these initial two lineages; otherwise  $t_1$  would not correspond to the  $t_{\text{MRCA}}$  of our study group. Each lineage evolves independently of all other lineages, giving rise to exactly one new lineage with rate  $b(t)$  and losing one existing lineage with rate  $d(t)$  (Figure 1 and Figure 2). Note that although each lineage evolves independently, all lineages share both a common (tree-wide) speciation rate  $b(t)$  and a common extinction rate  $d(t)$  ([Nee et al. 1994](#); [Höhna 2015](#)). Additionally, at certain times,  $t_{\text{M}}$ , a mass-extinction event occurs and each species existing at that time has the same probability,  $\rho$ , of survival. Finally, all extinct lineages are pruned and only the reconstructed tree remains (Figure 1).

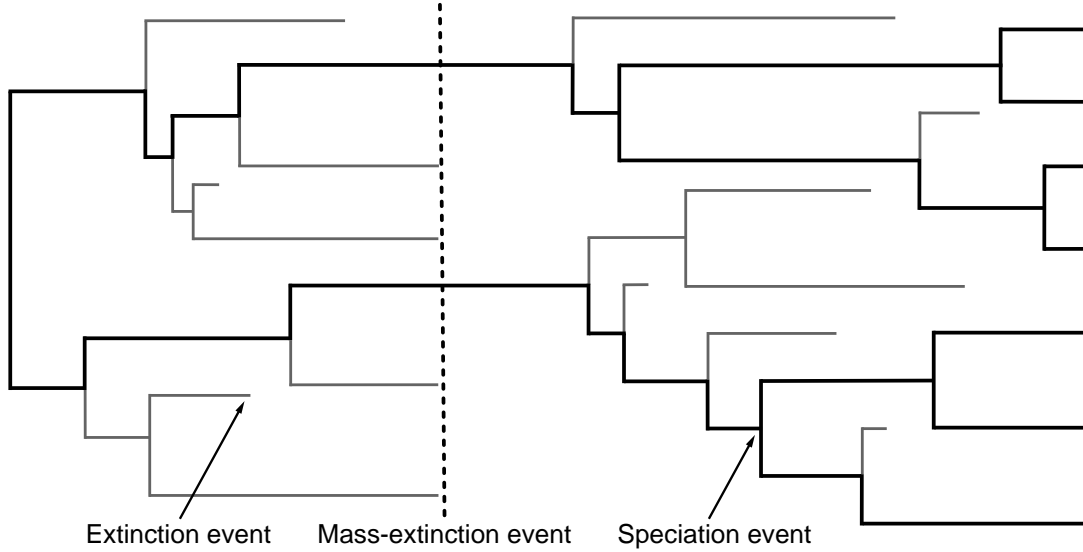


Figure 1: A realization of the birth-death process with mass extinction. Lineages that have no extant or sampled descendant are shown in gray and surviving lineages are shown in a thicker black line.

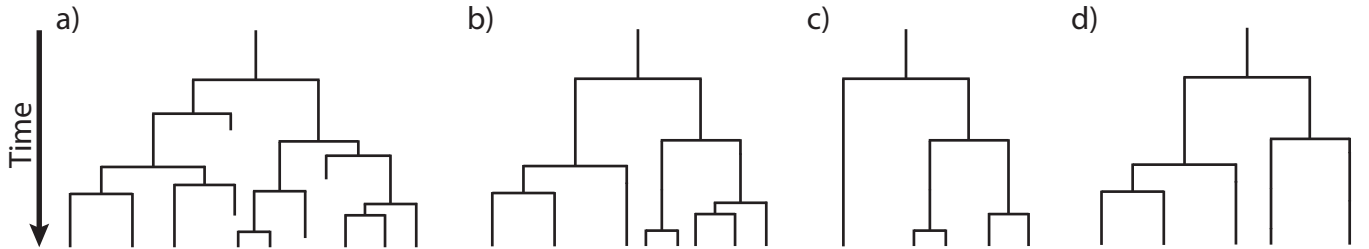


Figure 2: **Examples of trees produced under a birth-death process.** The process is initiated at the first speciation event (the ‘crown-age’ of the MRCA) when there are two initial lineages. At each speciation event the ancestral lineage is replaced by two descendant lineages. At an extinction event one lineage simply terminates. (A) A complete tree including extinct lineages. (B) The reconstructed tree of tree from A with extinct lineages pruned away. (C) A *uniform* subsample of the tree from B, where each species was sampled with equal probability,  $\rho$ . (D) A *diversified* subsample of the tree from B, where the species were selected so as to maximize diversity.

To condition the probability of observing the branching times on the survival of both lineages that descend from the root, we divide by  $P(N(T) > 0 | N(0) = 1)^2$ . Then, the probability density of the branching times,  $\mathbb{T}$ , becomes

$$P(\mathbb{T}) = \frac{\overbrace{P(N(T) = 1 \mid N(0) = 1)^2}^{\text{both initial lineages have one descendant}}}{\underbrace{P(N(T) > 0 \mid N(0) = 1)^2}_{\text{both initial lineages survive}}} \times \prod_{i=2}^{n-1} \overbrace{i \times b(t_i)}^{\text{speciation rate}} \times \overbrace{P(N(T) = 1 \mid N(t_i) = 1)}^{\text{lineage has one descendant}},$$

and the probability density of the reconstructed tree (topology and branching times) is then

$$P(\Psi) = \frac{2^{n-1}}{n!(n-1)!} \times \left( \frac{P(N(T) = 1 \mid N(0) = 1)}{P(N(T) > 0 \mid N(0) = 1)} \right)^2 \times \prod_{i=2}^{n-1} i \times b(t_i) \times P(N(T) = 1 \mid N(t_i) = 1) \quad (1)$$

We can expand Equation (1) by substituting  $P(N(T) > 0 \mid N(t) = 1)^2 \exp(r(t, T))$  for  $P(N(T) = 1 \mid N(t) = 1)$ , where  $r(u, v) = \int_u^v d(t) - b(t)dt$ ; the above equation becomes

$$\begin{aligned} P(\Psi) &= \frac{2^{n-1}}{n!(n-1)!} \times \left( \frac{P(N(T) > 0 \mid N(0) = 1)^2 \exp(r(0, T))}{P(N(T) > 0 \mid N(0) = 1)} \right)^2 \\ &\quad \times \prod_{i=2}^{n-1} i \times b(t_i) \times P(N(T) > 0 \mid N(t_i) = 1)^2 \exp(r(t_i, T)) \\ &= \frac{2^{n-1}}{n!} \times \left( P(N(T) > 0 \mid N(0) = 1) \exp(r(0, T)) \right)^2 \\ &\quad \times \prod_{i=2}^{n-1} b(t_i) \times P(N(T) > 0 \mid N(t_i) = 1)^2 \exp(r(t_i, T)). \end{aligned} \quad (2)$$

For a detailed description of this substitution, see [Höhna \(2015\)](#). Additional information regarding the underlying birth-death process can be found in ([Thompson 1975](#); Equation 3.4.6) and [Nee et al. \(1994\)](#) for constant rates and [Lambert \(2010\)](#); [Lambert and Stadler \(2013\)](#); [Höhna \(2013; 2014; 2015\)](#) for arbitrary rate functions.

To compute the equation above we need to know the rate function,  $r(t, s) = \int_t^s d(x) - b(x)dx$ , and the probability of survival,  $P(N(T) > 0 \mid N(t) = 1)$ . [Yule \(1925\)](#) and later [Kendall \(1948\)](#) derived the probability that a process survives ( $N(T) > 0$ ) and the probability of obtaining exactly  $n$  species at time  $T$  ( $N(T) = n$ ) when the process started at time  $t$  with one species. Kendall's results were summarized in Equation (3) and Equation (24) in [Nee et al. \(1994\)](#)

$$P(N(T) > 0 \mid N(t) = 1) = \left( 1 + \int_t^T \left( \mu(s) \exp(r(t, s)) \right) ds \right)^{-1} \quad (3)$$

$$\begin{aligned} P(N(T) = n \mid N(t) = 1) &= (1 - P(N(T) > 0 \mid N(t) = 1) \exp(r(t, T)))^{n-1} \\ &\quad \times P(N(T) > 0 \mid N(t) = 1)^2 \exp(r(t, T)) \end{aligned} \quad (4)$$

An overview for different diversification models is given in [Höhna \(2015\)](#).

### 3 Estimating Constant Speciation & Extinction Rates

#### 3.1 Outline

This tutorial describes how to specify basic branching-process models in **RevBayes**; two variants of the constant-rate birth-death process ([Yule 1925](#); [Kendall 1948](#); [Thompson 1975](#); [Nee et al. 1994](#); [Rannala and Yang 1996](#); [Yang and Rannala 1997](#); [Höhna 2015](#)). The probabilistic graphical model is given for each component of this tutorial. After each model is specified, you will estimate speciation and extinction rates using Markov chain Monte Carlo (MCMC). Finally, you will estimate the marginal likelihood of the model and evaluate the relative support using Bayes factors.

## 3.2 Requirements

We assume that you have read and hopefully completed the following tutorials:

- `RB_Getting_Started`
- `RB_Basics_Tutorial`
- `RB_BayesFactor_Tutorial`

Note that the `RB_Basics_Tutorial` introduces the basic syntax of `Rev` but does not cover any phylogenetic models. You may skip the `RB_Basics_Tutorial` if you have some familiarity with `R`. The `RB_BayesFactor_Tutorial` introduced Bayesian model selection by means of Bayes factors, which can be skipped by readers familiar with Bayesian model selection. We tried to keep this tutorial very basic and introduce all the language concepts and theory on the way. You may only need the `RB_Basics_Tutorial` for a more in-depth discussion of concepts in `Rev`.

## 4 Data and files

We provide the data file(s) which we will use in this tutorial. You may want to use your own data instead. In the `data` folder, you will find the following files

- `primates_springer.tre`: Dated primates phylogeny including 369 species from [Springer et al. \(2012\)](#).

→ Open the tree `data/primates_springer.tre` in FigTree.

## 5 Pure-Birth (Yule) Model

Before evaluating the relative support for different models, we must first specify them in `Rev`. In this section, we will walk through specifying a pure-birth process model and estimating the marginal likelihood. The section about the birth-death process will be less detailed because it will build up on this section.

The simplest branching model is the *pure-birth process* described by [Yule \(1925\)](#). Under this model, we assume at any instant in time, every lineage has the same speciation rate  $\lambda$ . In its simplest form, the speciation rate remains constant over time. As a result, the waiting time between speciation events is exponential, where the rate of the exponential distribution is the product of the number of extant lineages ( $n$ ) at that time and the speciation rate:  $n\lambda$  ([Yule 1925](#); [Aldous 2001](#); [Höhna 2014](#)). The pure-birth branching model does not allow for lineage extinction (*i.e.*, the extinction rate  $\mu = 0$ ). However, the model depends on a second parameter  $\rho$  which is the probability of sampling a species in the present time as well as the time of the start of the process, whether that is the origin time or root age. Therefore, the probabilistic graphical model of the pure-birth process is quite simple, where the observed time tree topology and node ages are conditional on the speciation rate, sampling probability, and root age (Fig. 3).

We can add hierarchical structure to this model and account for uncertainty in the value of the speciation rate by placing a hyperprior on  $\lambda$  (Fig. 4). The graphical models in Figures 3 and 4 demonstrate the

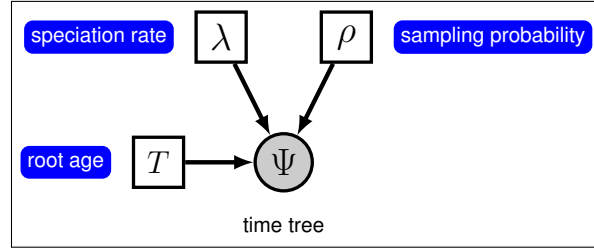


Figure 3: The graphical model representation of the pure-birth (Yule) process.

simplicity of the Yule model. Ultimately, the pure birth model is just a special case of the birth-death process, where the extinction rate (typically denoted  $\mu$ ) is a constant node with the value 0.

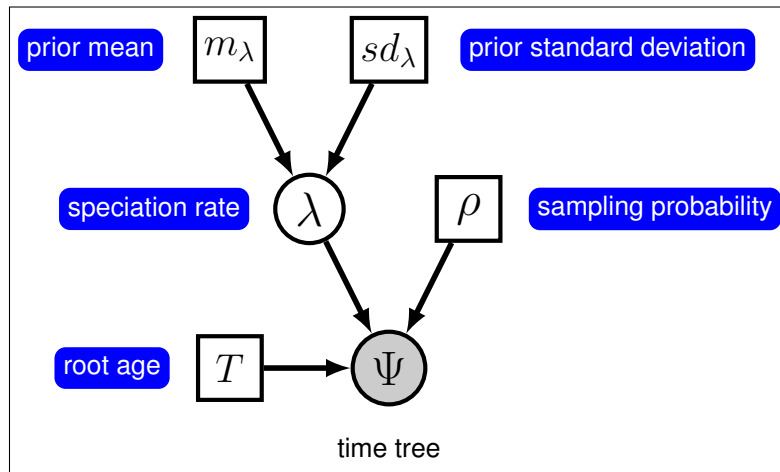


Figure 4: The graphical model representation of the pure-birth (Yule) process, where the speciation rate is treated as a random variable drawn from a lognormal distribution.

For this exercise, we will specify a Yule model, such that the speciation rate is a stochastic node, drawn from a lognormal distribution as in Figure 4. In a Bayesian framework, we are interested in estimating the posterior probability of  $\lambda$  given that we observe a time tree.

$$\mathbb{P}(\lambda \mid \Psi) = \frac{\mathbb{P}(\Psi \mid \lambda)\mathbb{P}(\lambda \mid \nu)}{\mathbb{P}(\Psi)} \quad (5)$$

In this example, we have a phylogeny of all living primates. We are treating the time tree  $\Psi$  as an observation, thus clamping the model with an observed value. The time tree we are conditioning the process on is taken from the analysis by [Springer et al. \(2012\)](#). Furthermore, there are approximately 450 described primates species, so we will fix the parameter  $\rho$  to 369/450.

- The full Yule-model specification is in the file called [Yule.Rev](#) on the [RevBayes](#) tutorial repository.

## 5.1 Read the tree

Begin by reading in the observed tree.

```
T <- readTrees("data/primates_springer.tre")[1]
```

From this tree, we can get some helpful variables:

```
taxa <- T.taxa()
```

Additionally, we can initialize an iterator variable for our vector of moves:

```
mvi = 0
mni = 0
```

## 5.2 Specifying the model

### 5.2.1 Birth rate

The model we are specifying only has three nodes (Fig. 4). We can specify the birth rate  $\lambda$ , the mean and standard deviation of the lognormal hyperprior on  $\lambda$ , and the conditional dependency of the two parameters all in one line of Rev code.

```
birth_rate_mean <- ln( ln(450/2) / T.rootAge() )
birth_rate_sd <- 0.587405
birth_rate ~ dnLognormal(mean=birth_rate_mean,sd=birth_rate_sd)
```

Here, the stochastic node called **birth\_rate** represents the speciation rate  $\lambda$ . **birth\_rate\_mean** and **birth\_rate\_sd** are the prior mean and prior standard deviation, respectively. We chose the prior mean so that it is centered around observed number of species (*i.e.*, the expected number of species under a Yule process will thus be equal to the observed number of species) and a prior standard deviation of 0.587405 which creates a lognormal distribution with 95% prior probability spanning exactly one magnitude. If you want to represent more prior uncertainty by, *e.g.*, allowing for two orders of magnitude in the 95% prior probability then you can simply multiply **birth\_rate\_sd** by a factor of 2.

To estimate the value of  $\lambda$ , we assign a proposal mechanism to operate on this node. In RevBayes these MCMC sampling algorithms are called *moves*. We need to create a vector of moves and we can do this by using vector indexing and our pre-initialized iterator **mi**. We will use a scaling move on  $\lambda$  called **mvScale**.

```
moves[++mvi] = mvScale(birth_rate,lambda=1,tune=true,weight=3)
```

### 5.2.2 Sampling probability

Our prior belief is that we have sampled 367 out of 450 living primate species. To account for this we can set the sampling parameter as a constant node with a value of 369/450

```
rho <- T.ntips()/450
```

### 5.2.3 Root age

Any stochastic branching process must be conditioned on a time that represents the start of the process. Typically, this parameter is the *origin time* and it is assumed that the process started with *one* lineage. Thus, the origin of a birth-death process is the node that is *ancestral* to the root node of the tree. For macroevolutionary data, particularly without any sampled fossils, it is difficult to use the origin time. To accommodate this, we can condition on the age of the root by assuming the process started with *two* lineages that both originate at the time of the root.

We can get the value for the root from the [Springer et al. \(2012\)](#) tree.

```
root_time <- T.rootAge()
```

### 5.2.4 The time tree

Now we have all of the parameters we need to specify the full pure-birth model. We can initialize the stochastic node representing the time tree. Note that we set the **mu** parameter to the constant value **0.0**.

```
timetree ~ dnBDP(lambda=birth_rate, mu=0.0, rho=rho, rootAge=root_time,
  samplingStrategy="uniform", condition="survival", taxa=taxa)
```

If you refer back to Equation 5 and Figure 4, the time tree  $\Psi$  is the variable we observe, *i.e.*, the data. We can set this in **Rev** by using the **clamp()** function.

```
timetree.clamp(T)
```

Here we are fixing the value of the time tree to our observed tree from [Springer et al. \(2012\)](#).

Finally, we can create a workspace object of our whole model using the **model()** function. Workspace objects are initialized using the **=** operator. This distinguishes the objects used by the program to run the MCMC analysis from the distinct nodes of our graphical model. The model workspace objects makes it easy to work with the model in **Rev** and creates a wrapper around our model DAG. Because our model is a directed, acyclic graph (DAG), we only need to give the model wrapper function a single node and it does the work to find all the other nodes through their connections.

```
mymodel = model(birth_rate)
```

The **model()** function traversed all of the connections and found all of the nodes we specified.



## 5.3 Running an MCMC analysis

### 5.3.1 Specifying Monitors

For our MCMC analysis, we need to set up a vector of *monitors* to record the states of our Markov chain. The monitor functions are all called **mn\***, where **\*** is the wildcard representing the monitor type. First, we will initialize the model monitor using the **mnModel** function. This creates a new monitor variable that will output the states for all model parameters when passed into a MCMC function.

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

Additionally, create a screen monitor that will report the states of specified variables to the screen with **mnScreen**:

```
monitors[++mni] = mnScreen(printgen=1000, birth_rate)
```

### 5.3.2 Initializing and Running the MCMC Simulation

With a fully specified model, a set of monitors, and a set of moves, we can now set up the MCMC algorithm that will sample parameter values in proportion to their posterior probability. The **mcmc()** function will create our MCMC object:

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

We may wish to run the **.burnin()** member function, *i.e.*, if we wish to pre-run the chain and discard the initial states. Recall that the **.burnin()** function specifies a *completely separate* preliminary MCMC analysis that is used to tune the scale of the moves to improve mixing of the MCMC analysis.

```
mymcmc.burnin(generations=10000, tuningInterval=200)
```

Now, run the MCMC:

```
mymcmc.run(generations=50000)
```

When the analysis is complete, you will have the monitored files in your output directory.

→ The Rev file for performing this analysis: [mcmc\\_Yule.Rev](#).

## 5.4 Exercise 1

- Run an MCMC simulation to estimate the posterior distribution of the speciation rate (**birth\_rate**).
- Load the generated output file into **Tracer**: What is the mean posterior estimate of the **birth\_rate** and what is the estimated HPD?
- Compare the prior mean with the posterior mean. (**Hint**: Use the optional argument **underPrior=TRUE** in the function **mymcmc.run()**) Are they different (*e.g.*, Figure 5)? Is the posterior mean outside the prior 95% probability interval?
- Repeat the analysis and allow for two orders of magnitude of prior uncertainty.

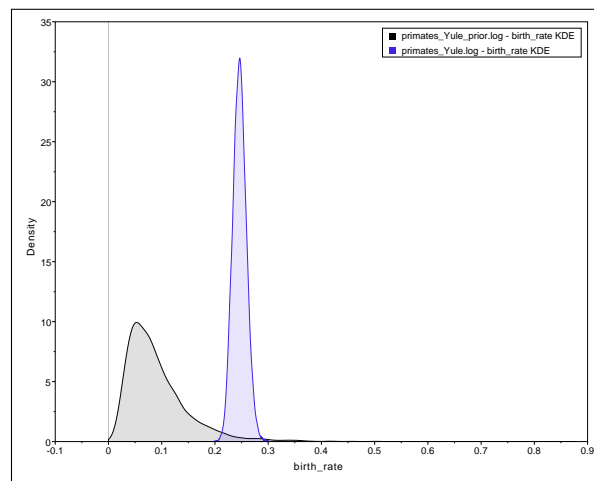


Figure 5: Estimates of the posterior and prior distribution of the **birth\_rate** visualized in **Tracer**. The prior (black curve) shows lognormal distribution that we chose as the prior distribution.

## 6 Estimating the marginal likelihood of the model

With a fully specified model, we can set up the **powerPosterior()** analysis to create a file of ‘powers’ and likelihoods from which we can estimate the marginal likelihood using stepping-stone or path sampling. This method computes a vector of powers from a beta distribution, then executes an MCMC run for each power step while raising the likelihood to that power. In this implementation, the vector of powers starts with 1, sampling the likelihood close to the posterior and incrementally sampling closer and closer to the prior as the power decreases. For more information on marginal likelihood estimation please read the [RB\\_BayesFactor\\_Tutorial](#).

First, we create the variable containing the power posterior. This requires us to provide a model and vector of moves, as well as an output file name. The **cats** argument sets the number of power steps.

```
pow_p = powerPosterior(mymodel, moves, monitors, "output/Yule_powp.out", cats=100,
    sampleFreq=10)
```

We can start the power posterior by first burning in the chain and discarding the first 10000 states.

```
pow_p.burnin(generations=10000,tuningInterval=200)
```

Now execute the run with the `.run()` function:

```
pow_p.run(generations=10000)
```

Once the power posteriors have been saved to file, create a stepping stone sampler. This function can read any file of power posteriors and compute the marginal likelihood using stepping-stone sampling.

```
ss = steppingStoneSampler(file="output/Yule_powp.out", powerColumnName="power",
    likelihoodColumnName="likelihood")
```

Compute the marginal likelihood under stepping-stone sampling using the member function `marginal()` of the `ss` variable and record the value in Table 1.

```
ss.marginal()
```

Path sampling is an alternative to stepping-stone sampling and also takes the same power posteriors as input.

```
ps = pathSampler(file="output/Yule_powp.out", powerColumnName="power", likelihoodColumnName="
    likelihood")
```

Compute the marginal likelihood under stepping-stone sampling using the member function `marginal()` of the `ps` variable and record the value in Table 1.

```
ps.marginal()
```

→ The Rev file for performing this analysis: [ml\\_Yule.Rev](#).

## 6.1 Exercise 2

- Compute the marginal likelihood under the Yule model.
- Enter the estimate in the table below.

Table 1: Marginal likelihoods and Bayes factors\*.

Estimate	Stepping-stone	Path sampling
Marginal likelihood Yule ( $M_0$ )		
Marginal likelihood birth-death ( $M_1$ )		
Supported model?		

\*you can edit this table

## 7 Birth-Death Process

The pure-birth model does not account for extinction, thus it assumes that every lineage at the start of the process will have sampled descendants at time 0. This assumption is fairly unrealistic for most phylogenetic datasets on a macroevolutionary time scale since the fossil record provides evidence of extinct lineages. [Kendall \(1948\)](#) described a more general branching process model to account for lineage extinction called the *birth-death process*. Under this model, at any instant in time, every lineage has the same rate of speciation  $\lambda$  and the same rate of extinction  $\mu$ . This is the *constant-rate* birth-death process, which considers the rates constant over time and over the tree ([Nee et al. 1994](#); [Höhna 2015](#)).

[Yang and Rannala \(1997\)](#) derived the probability of time trees under an extension of the birth-death model that accounts for incomplete sampling of the tips (Fig. 6) (see also [Stadler \(2009\)](#) and [Höhna \(2014\)](#)). Under this model, the parameter  $\rho$  accounts for the probability of sampling in the present time, and because it is a probability, this parameter can only take values between 0 and 1.

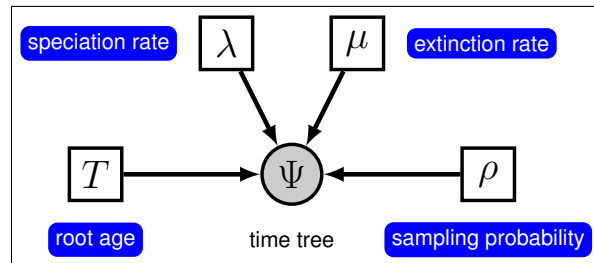


Figure 6: The graphical model representation of the birth-death process with uniform sampling and conditioned on the root age.

In principle, we can specify a model with prior distributions on speciation and extinction rates directly. One possibility is to specify an exponential, lognormal, or gamma distribution as the prior on either rate parameter. However, it is more common to specify prior distributions on a transformation of the speciation and extinction rate because, for example, we want to enforce that the speciation rate is always larger than the extinction rate.

In the following subsections we will only provide the key command that are different for the constant-rate birth-death process. All other commands will be the same as in the previous exercise. You should copy the `mcmc_Yule.Rev` script and modify it accordingly. Don't forget to rename the filenames of the monitors to avoid overwriting of your previous results!

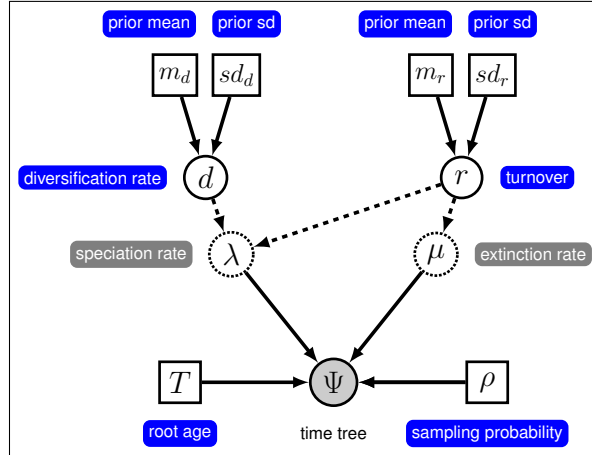


Figure 7: The graphical model representation of the birth-death process with uniform sampling parameterized using the diversification and turnover.

## 7.1 Diversification and turnover

We have some good prior information about the magnitude of the diversification. The diversification rate represent the rate at which the species diversity increases. Thus, we just use the same prior for the diversification rate as we used before for the birth rate.

```
diversification_mean <- ln( ln(450.0/2.0) / T.rootAge() )
diversification_sd <- 0.587405
diversification ~ dnLognormal(mean=diversification_mean,sd=diversification_sd)
moves[+mvi] = mvScale(diversification,lambda=1.0,tune=true,weight=3.0)
```

Unfortunately, we have less prior information about the turnover rate. The turnover rate is the rate how fast one species is replaced by another species due to a birth plus death event. Hence, the turnover rate represent the longevity of a species. For simplicity we use the same prior on the turnover rate but with two orders of magnitude prior uncertainty.

```
turnover_mean <- ln( ln(450.0/2.0) / T.rootAge() )
turnover_sd <- 0.587405*2
turnover ~ dnLognormal(mean=turnover_mean,sd=turnover_sd)
moves[+mvi] = mvScale(turnover,lambda=1.0,tune=true,weight=3.0)
```

## 7.2 Birth rate and death rate

The birth and death rates are both deterministic nodes. We compute them by simple parameter transformation. Note that the death rate is in fact equal to the turnover rate.

```
birth_rate := diversification + turnover
death_rate := turnover
```

All other parameters, such as the sampling probability and the root age are kept the same as in the analysis above.

### 7.3 The time tree

Initialize the stochastic node representing the time tree. The main difference now is that we provide a stochastic parameter for the extinction rate  $\mu$ .

```
timetree ~ dnBDP(lambda=birth_rate, mu=death_rate, rho=rho, rootAge=root_time,
  samplingStrategy="uniform", condition="survival", taxa=taxa)
```

### 7.4 Exercise 3

- Run an MCMC simulation to compute the posterior distribution of the diversification and turnover rate.
- Look at the parameter estimates in **Tracer**. What can you say about the diversification, turnover, speciation and extinction rates? How high is the extinction rate compared with the speciation rate?
- Compute the marginal likelihood under the BD model. Which model is supported by the data?
- Enter the estimate in the table above.
- Can you modify the script to use a prior on the birth drawn from a lognormal distribution and relative death rate drawn from a beta distribution so that the extinction rate is equal to the birth rate times the relative death rate?
  - a) Do the parameter estimates change?
  - b) What about the marginal likelihood estimates?

## References

- Aldous, D. J. 2001. Stochastic models and descriptive statistics for phylogenetic trees, from Yule to today. *Statistical Science* Pages 23–34.
- Drummond, A., M. Suchard, D. Xie, and A. Rambaut. 2012. Bayesian phylogenetics with *beast* and the *beast* 1.7. *Molecular Biology and Evolution* 29:1969–1973.
- Höhna, S. 2013. Fast simulation of reconstructed phylogenies under global time-dependent birth-death processes. *Bioinformatics* 29:1367–1374.
- Höhna, S. 2014. Likelihood Inference of Non-Constant Diversification Rates with Incomplete Taxon Sampling. *PLoS One* 9:e84184.

- Höhna, S. 2015. The time-dependent reconstructed evolutionary process with a key-role for mass-extinction events. *Journal of Theoretical Biology* 380:321–331.
- Höhna, S., M. J. Landis, B. Boussau, B. R. Moore, N. Lartillot, T. A. Heath, J. P. Huelsenbeck, and F. Ronquist. 2015. RevBayes: Bayesian Phylogenetic Inference Using Graphical Models and an Interactive Model Specification Language. submitted .
- Kendall, D. G. 1948. On the generalized "birth-and-death" process. *The Annals of Mathematical Statistics* 19:1–15.
- Lambert, A. 2010. The contour of splitting trees is a lévy process. *The Annals of Probability* 38:348–395.
- Lambert, A. and T. Stadler. 2013. Birth–death models and coalescent point processes: the shape and probability of reconstructed phylogenies. *Theoretical Population Biology* 90:113–128.
- Nee, S., R. M. May, and P. H. Harvey. 1994. The Reconstructed Evolutionary Process. *Philosophical Transactions: Biological Sciences* 344:305–311.
- Rannala, B. and Z. Yang. 1996. Probability distribution of molecular evolutionary trees: A new method of phylogenetic inference. *Journal of Molecular Evolution* 43:304–311.
- Ronquist, F., M. Teslenko, P. van der Mark, D. L. Ayres, A. Darling, S. Höhna, B. Larget, L. Liu, M. A. Suchard, and J. P. Huelsenbeck. 2012. Mrbayes 3.2: efficient bayesian phylogenetic inference and model choice across a large model space. *Systematic Biology* 61:539–542.
- Springer, M. S., R. W. Meredith, J. Gatesy, C. A. Emerling, J. Park, D. L. Rabosky, T. Stadler, C. Steiner, O. A. Ryder, J. E. Janečka, et al. 2012. Macroevolutionary dynamics and historical biogeography of primate diversification inferred from a species supermatrix. *PLoS One* 7:e49521.
- Stadler, T. 2009. On incomplete sampling under birth-death models and connections to the sampling-based coalescent. *Journal of Theoretical Biology* 261:58–66.
- Thompson, E. 1975. Human evolutionary trees. Cambridge University Press Cambridge.
- Yang, Z. and B. Rannala. 1997. Bayesian phylogenetic inference using DNA sequences: a Markov Chain Monte Carlo Method. *Molecular Biology and Evolution* 14:717–724.
- Yule, G. 1925. A mathematical theory of evolution, based on the conclusions of dr. jc willis, frs. *Philosophical Transactions of the Royal Society of London. Series B, Containing Papers of a Biological Character* 213:21–87.

Version dated: July 10, 2016