

Phylogenetic Inference using RevBayes

Model Selection & Data Partitioning

Overview

This tutorial demonstrates how to set up and perform an analysis that calculates Bayes factors to select among different partition configurations of aligned DNA sequences. After selecting the model that is best supported by the data, the exercise continues with basic inference of an unrooted tree topology and branch lengths using Markov chain Monte Carlo (MCMC).

Requirements

We assume that you have completed the following tutorials:

- [RB_Basics_Tutorial](#)
- [RB_CTMC_Tutorial](#)
- [RB_MCMC_Tutorial](#)

1 Exercise: Model Selection & Partitioning using Bayes Factors

1.1 Introduction

Variation in the evolutionary process across the sites of nucleotide sequence alignments is well established, and is an increasingly pervasive feature of datasets composed of gene regions sampled from multiple loci and/or different genomes. Inference of phylogeny from these data demands that we adequately model the underlying process heterogeneity; failure to do so can lead to biased estimates of phylogeny and other parameters (Brown and Lemmon, 2007). To accommodate process heterogeneity within and/or between various gene(omic) regions, we will evaluate the support for various partition schemes using Bayes factors to compare the marginal likelihoods of the candidate partition schemes.

Accounting for process heterogeneity involves adopting a ‘mixed-model’ approach, (Ronquist and Huelsenbeck, 2003) in which the sequence alignment is first parsed into a number of partitions that are intended to capture plausible process heterogeneity within the data. The determination of the partitioning scheme is guided by biological considerations regarding the dataset at hand. For example, we might wish to evaluate possible variation in the evolutionary process within a single gene region (*e.g.*, between stem and loop regions of ribosomal sequences), or among gene regions in a concatenated alignment (*e.g.*, comprising multiple nuclear loci and/or gene regions sampled from different genomes). The choice of partitioning scheme is up to the investigator and many possible partitions might be considered for a typical dataset.

Next, a substitution model is specified for each predefined process partition (using a given model-selection criterion, such as Bayes factors). In this exercise, we assume that each partition evolved under an independent general-time reversible model with gamma-distributed rates across sites (GTR+ Γ). Under this model the observed data are conditionally dependent on the exchangeability rates (θ), stationary base frequencies (π), and the degree of gamma-distributed among-site rate variation (α), as well as the unrooted tree topology (Ψ) and branch lengths (ν). We show the graphical model representation of the GTR+ Γ mode in Figure ???. When we assume different GTR+ Γ models for each partitions, this results in a composite model, in which all sites are assumed to share a common, unrooted tree topology and proportional branch lengths, but subsets of sites (‘data partitions’) are assumed to have independent substitution model parameters. This composite model is referred to as a *mixed model*.

Finally, we perform a separate MCMC simulation to approximate the joint posterior probability density of the phylogeny and other parameters. Note that, in this approach, the mixed model is a fixed assumption of the inference (*i.e.*, the parameter estimates are conditioned on the specified mixed model), and the parameters for each process partition are independently estimated.

For most sequence alignments, several (possibly many) partition schemes of varying complexity are plausible *a priori*, which therefore requires a way to objectively identify the partition scheme that balances estimation bias and error variance associated with under- and over-parameterized mixed models, respectively. Increasingly, mixed-model selection is based on *Bayes factors* (*e.g.*, Suchard, Weiss and Sinsheimer, 2001), which involves first calculating the marginal likelihood under each candidate partition scheme and then comparing the ratio of the marginal likelihoods for the set of candidate partition schemes (Brandley, Schmitz and Reeder, 2005; Nylander et al., 2004; McGuire et al., 2007). The analysis pipeline that we will use in this tutorial is depicted in Figure 1.

Given two models, M_0 and M_1 , the Bayes factor comparison assessing the relative plausibility of each

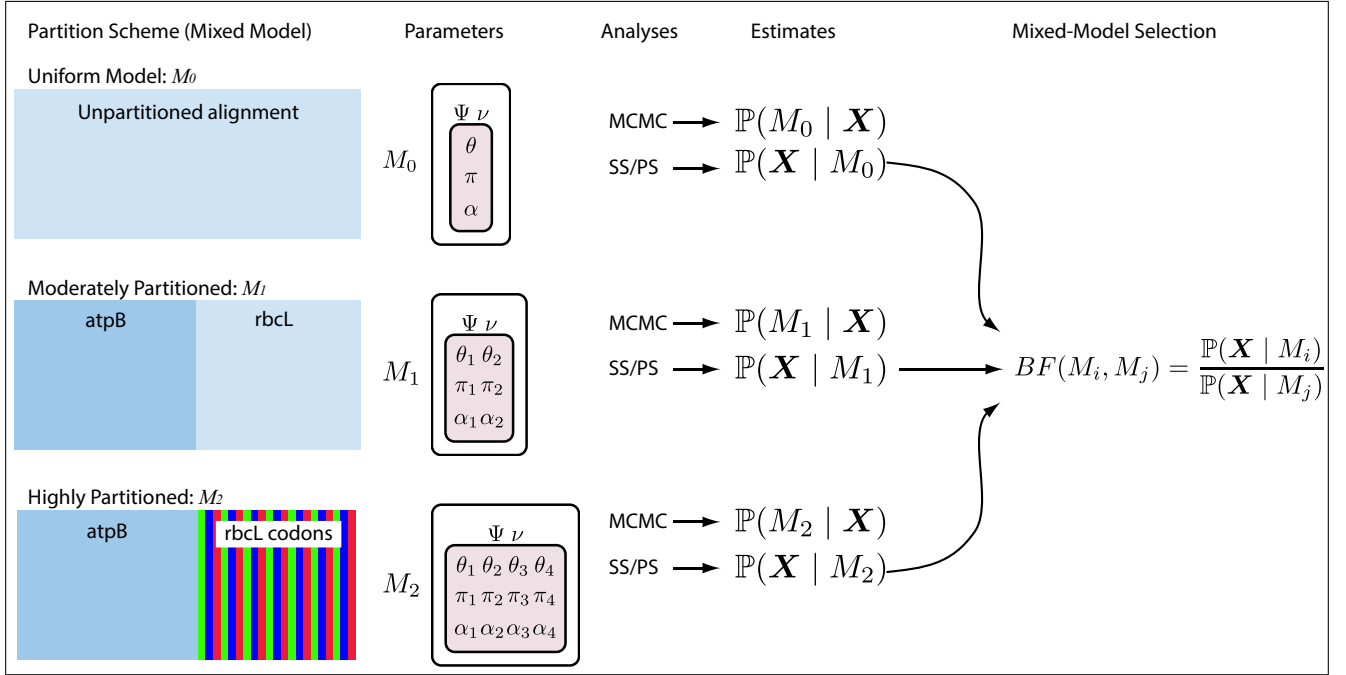


Figure 1: The analysis pipeline for Exercise 1. We will explore three partition schemes for the conifer dataset. The first model (the ‘uniform model’, M_0) assumes that all sites evolved under a common GTR+ Γ substitution model. The second model (the ‘moderately partitioned’ model, M_1) invokes two data partitions corresponding to the two gene regions (atpB and rbcL), and assumes each subset of sites evolved under an independent GTR+ Γ model. The final mixed model (the ‘highly partitioned’ model, M_2) invokes four data partitions—the first partition corresponds to the atpB gene region, and the remaining partitions correspond to the three codon positions of the rbcL gene region—and each data partition is assumed evolved under an independent GTR+ Γ substitution model. Note that we assume that all sites share a common tree topology, Ψ , and branch-length proportions, ν , for each of the candidate partition schemes. We perform two separate sets of analyses for each mixed model—a Metropolis-coupled MCMC simulation to approximate the joint posterior probability density of the mixed-model parameters, and a ‘stepping-stone’ MCMC simulation to approximate the marginal likelihood for each mixed model. The resulting marginal-likelihood estimates are then evaluated using Bayes factors to assess the fit of the data to the three candidate mixed models.

model as an explanation of the data, $BF(M_0, M_1)$, is:

$$BF(M_0, M_1) = \frac{\text{posterior odds}}{\text{prior odds}}.$$

The posterior odds is the posterior probability of M_0 given the data, \mathbf{X} , divided by the posterior odds of M_1 given the data:

$$\text{posterior odds} = \frac{\mathbb{P}(M_0 | \mathbf{X})}{\mathbb{P}(M_1 | \mathbf{X})},$$

and the prior odds is the prior probability of M_0 divided by the prior probability of M_1 :

$$\text{prior odds} = \frac{\mathbb{P}(M_0)}{\mathbb{P}(M_1)}.$$

Thus, the Bayes factor measures the degree to which the data alter our belief regarding the support for M_0 relative to M_1 (Lavine and Schervish, 1999):

$$BF(M_0, M_1) = \frac{\mathbb{P}(M_0 | \mathbf{X}, \theta_0)}{\mathbb{P}(M_1 | \mathbf{X}, \theta_1)} \div \frac{\mathbb{P}(M_0)}{\mathbb{P}(M_1)}. \quad (1)$$

This, somewhat vague, definition does not lead to clear-cut identification of the “best” model. Instead, you must decide the degree of your belief in M_0 relative to M_1 . Despite the absence of any strict “rule-of-thumb”, you can refer to the scale (outlined by [Jeffreys, 1961](#)) for interpreting these measures (Table 1).

Table 1: The scale for interpreting Bayes factors by Harold [Jeffreys \(1961\)](#).

$BF(M_0, M_1)$	Strength of evidence
$< 1 : 1$	Negative (supports M_1)
$1 : 1$ to $3 : 1$	Barely worth mentioning
$3 : 1$ to $10 : 1$	Substantial
$10 : 1$ to $30 : 1$	Strong
$30 : 1$ to $100 : 1$	Very strong
$> 100 : 1$	Decisive

For a detailed description of Bayes factors see [Kass and Raftery \(1995\)](#)

Unfortunately, direct calculation of the posterior odds to prior odds ratio is unfeasible for most phylogenetic models. However, we can further define the posterior odds ratio as:

$$\frac{\mathbb{P}(M_0 | \mathbf{X})}{\mathbb{P}(M_1 | \mathbf{X})} = \frac{\mathbb{P}(M_0) \mathbb{P}(\mathbf{X} | M_0)}{\mathbb{P}(M_1) \mathbb{P}(\mathbf{X} | M_1)},$$

where $\mathbb{P}(\mathbf{X} | M_i)$ is the *marginal likelihood* of the data marginalized over all parameters for M_i ; it is also referred to as the *model evidence* or *integrated likelihood*. More explicitly, the marginal likelihood is the probability of the set of observed data (\mathbf{X}) under a given model (M_i), while averaging over all possible values of the parameters of the model (θ_i) with respect to the prior density on θ_i

$$\mathbb{P}(\mathbf{X} | M_i) = \int \mathbb{P}(\mathbf{X} | \theta_i) \mathbb{P}(\theta_i) d\theta_i. \quad (2)$$

If you refer back to equation 1, you can see that, with very little algebra, the ratio of marginal likelihoods is equal to the Bayes factor:

$$BF(M_0, M_1) = \frac{\mathbb{P}(\mathbf{X} | M_0)}{\mathbb{P}(\mathbf{X} | M_1)} = \frac{\mathbb{P}(M_0 | \mathbf{X}, \theta_0)}{\mathbb{P}(M_1 | \mathbf{X}, \theta_1)} \div \frac{\mathbb{P}(M_0)}{\mathbb{P}(M_1)}. \quad (3)$$

Therefore, we can perform a Bayes factor comparison of two models by calculating the marginal likelihood for each one. Alas, exact solutions for calculating marginal likelihoods are not known for phylogenetic models (see equation 2), thus we must resort to numerical integration methods to estimate or approximate these values. In this exercise, we will estimate the marginal likelihood for each partition scheme using both the stepping-stone ([Xie et al., 2011](#)) and path sampling estimators ([Gelman and Meng, 1998](#); [Lartillot and Philippe, 2006](#); [Friel and Pettitt, 2008](#)).

1.2 Getting Started

The various exercises in this tutorial take you through the steps required to perform phylogenetic analyses of the example datasets. In addition, we have provided the output files for every exercise so you can verify your results. (Note that since the MCMC runs you perform will start from different random seeds, the output files resulting from your analyses *will not* be identical to the ones we provide you.)

- Download data and output files from:
https://www.nescent.org/sites/academy/RevBayes_Workshop_Schedule

- Also note that “pre-cooked” output files are provided in the download. Throughout this tutorial, you can use those files to summarize output if you do not have time to run the full analyses yourself.

1.3 Launch RevBayes

Execute the RevBayes binary. If this program is in your path, then you can simply type in your Unix terminal:

- `$ rb`

1.4 Phylogenetic Models

The models we use here are equivalent to the models described in the previous exercise on MCMC methods and convergence assessment. To specify the model please consult the previous exercise.

ESTIMATING THE MARGINAL LIKELIHOOD

Typically, model comparison is performed prior to running the full MCMC analysis under a model. If you calculated the Bayes factors to determine the relative support for the uniform model and found that there was strong evidence supporting this model over others (hint: this is not true if you proceed with this tutorial), then it would be worth your time to proceed with the MCMC steps outlined above. The following steps will describe using stepping-stone and path sampling methods on a set of power posteriors to estimate marginal likelihoods under the uniform 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.

Just to be safe, it is better to clear the workspace and re-load the data and model:

```
RevBayes > clear()
RevBayes > source("RevBayes_scripts/uniform_partition_model.Rev")
```

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.

```
RevBayes > pow_p <- powerPosterior(mymodel, moves, "pow_p_PS0.out", cats=50)
```

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

```
RevBayes > pow_p.burnin(generations=10000,tuningInterval=1000)
```

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

```
RevBayes > pow_p.run(generations=1000)
```

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.

```
RevBayes > ss <- steppingStoneSampler(file="pow_p_PS0.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 2.

```
RevBayes > ss.marginal()
```

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

```
RevBayes > ps <- pathSampler(file="pow_p_PS0.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 2.

```
RevBayes > ps.marginal()
```

Repeat this analysis for each of the other three partitioning schemes.

1.5 Compute Bayes Factors and Select Model

Now that we have estimates of the marginal likelihood under each of our different models, we can evaluate their relative plausibility using Bayes factors. Use Table 2 to summarize the marginal log-likelihoods estimated using the stepping-stone and path-sampling methods.

Table 2: Estimated marginal likelihoods for different partition configurations*.

Partition	Marginal lnL estimates	
	<i>Stepping-stone</i>	<i>Path sampling</i>
1.4 uniform (M_1)		
?? moderate (M_2)		
?? extreme (M_3)		

*you can edit this table

Phylogenetics software programs log-transform the likelihood to avoid [underflow](#), because multiplying likelihoods results in numbers that are too small to be held in computer memory. Thus, we must use a different form of equation 3 to calculate the ln-Bayes factor (we will denote this value \mathcal{K}):

$$\mathcal{K} = \ln[BF(M_0, M_1)] = \ln[\mathbb{P}(\mathbf{X} \mid M_0)] - \ln[\mathbb{P}(\mathbf{X} \mid M_1)], \quad (4)$$

where $\ln[\mathbb{P}(\mathbf{X} \mid M_0)]$ is the *marginal lnL* estimate for model M_0 . The value resulting from equation 4 can be converted to a raw Bayes factor by simply taking the exponent of \mathcal{K}

$$BF(M_0, M_1) = e^{\mathcal{K}}. \quad (5)$$

Alternatively, you can interpret the strength of evidence in favor of M_0 using the \mathcal{K} and skip equation 5. In this case, we evaluate the \mathcal{K} in favor of model M_0 against model M_1 so that:

$$\begin{aligned} &\text{if } \mathcal{K} > 1, \text{ then model } M_0 \text{ wins} \\ &\text{if } \mathcal{K} < -1, \text{ then model } M_1 \text{ wins.} \end{aligned}$$

Thus, values of \mathcal{K} around 0 indicate ambiguous support.

Using the values you entered in Table 2 and equation 4, calculate the ln-Bayes factors (using \mathcal{K}) for the different model comparisons. Enter your answers in Table 3 using the stepping-stone and the path-sampling estimates of the marginal log likelihoods.

Because of the computational costs of computing marginal likelihoods and the vast number of possible partitioning strategies, it is not feasible to evaluate all of them. New methods based on nonparametric Bayesian models have recently been applied to address this problem ([Lartillot and Philippe, 2004](#); [Huelsenbeck and Suchard, 2007](#); [Wu, Suchard and Drummond, 2013](#)). These approaches use an infinite mixture model (the Dirichlet process; [Ferguson, 1973](#); [Antoniak, 1974](#)) that places non-zero probability on *all* of the countably-infinite possible partitions for a set of sequences. Bayesian phylogenetic inference under these models is implemented in the program [PhyloBayes](#) ([Lartillot, Lepage and Blanquart, 2009](#)) and the [subst-bma](#) plug-in for [BEAST2](#) ([Wu, Suchard and Drummond, 2013](#)).

Note that Bayes factors based on comparison of HM-based marginal likelihoods often *strongly* favor the most extremely partitioned mixed model. In fact, the harmonic mean estimator has been shown to provide unreliable estimates of marginal likelihoods, compared to more robust approaches ([Lartillot and Philippe, 2006](#); [Xie et al., 2011](#); [Fan et al., 2011](#)). Based on these studies, it is recommended that you avoid using HM-derived marginal likelihoods for Bayes factor comparisons. (The Canadian Bayesian Radford Neal says the harmonic mean is the “[worst Monte Carlo method ever](#)”.)

Table 3: Bayes factor calculation*.

Model comparison	ln-Bayes Factor (\mathcal{K})	
	<i>Stepping-stone</i>	<i>Path sampling</i>
M_1, M_2		
M_1, M_3		
M_1, M_4		
M_2, M_3		
M_2, M_4		
M_3, M_4		
Supported model?		

*you can edit this table

Batch Mode

If you wish to run this exercise in batch mode, the files are provided for you.

You can carry out these batch commands by providing the file name when you execute the **rb** binary in your unix terminal (this will overwrite all of your existing run files).


- `$ rb full_analysis.Rev`

Useful Links

- RevBayes: <https://github.com/revbayes/code>
- MrBayes: <http://mrbayes.sourceforge.net>
- PhyloBayes: <http://www.phylobayes.org>
- Bali-Phy: <http://www.bali-phy.org>
- Tree Thinkers: <http://treethinkers.org>

Questions about this tutorial can be directed to:

- Tracy Heath (email: tracyh@berkeley.edu)
- Michael Landis (email: mlandis@berkeley.edu)
- Sebastian Höhna (email: sebastian.hohna@gmail.com)

 This tutorial was written by [Tracy Heath](#), [Michael Landis](#), and Sebastian Höhna; licensed under a [Creative Commons Attribution 4.0 International License](#). (This tutorial is based on the [Phylogenetic Inference using MrBayes v3.2](#) tutorial written by Tracy Heath, Conor Meehan, and Brian Moore and some content is reproduced here. Mark Holder, Ben Liebeskind, Emily McTavish, and April Wright provided helpful comments.)

Version dated: August 26, 2014

Relevant References

- Antoniak CE. 1974. Mixtures of Dirichlet processes with applications to non-parametric problems. *Annals of Statistics*. 2:1152–1174.
- Brandley MC, Schmitz A, Reeder TW. 2005. Partitioned Bayesian analyses, partition choice, and the phylogenetic relationships of scincid lizards. *Systematic Biology*. 54:373–390.
- Brown JM, Lemmon AR. 2007. The importance of data partitioning and the utility of Bayes factors in Bayesian phylogenetics. *Systematic Biology*. 56:643–655.
- Fan Y, Wu R, Chen MH, Kuo L, Lewis PO. 2011. Choosing among partition models in Bayesian phylogenetics. *Molecular Biology and Evolution*. 28:523–532.
- Ferguson TS. 1973. A Bayesian analysis of some nonparametric problems. *Annals of Statistics*. 1:209–230.
- Friel N, Pettitt AN. 2008. Marginal likelihood estimation via power posteriors. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*. 70:589–607.
- Gelman A, Meng XL. 1998. Simulating normalizing constants: From importance sampling to bridge sampling to path sampling. *Statistical science*. pp. 163–185.
- Huelsenbeck JP, Suchard M. 2007. A nonparametric method for accommodating and testing across-site rate variation. *Systematic Biology*. 56:975–987.
- Jeffreys H. 1961. *Theory of Probability*. Oxford: Oxford University Press.
- Kass RE, Raftery AE. 1995. Bayes factors. *Journal of the American Statistical Association*. 90:773–795.
- Lartillot N, Lepage T, Blanquart S. 2009. Phylobayes 3: a bayesian software package for phylogenetic reconstruction and molecular dating. *Bioinformatics*. 25:2286.
- Lartillot N, Philippe H. 2004. A Bayesian mixture model for across-site heterogeneities in the amino-acid replacement process. *Molecular Biology and Evolution*. 21:1095–1109.
- Lartillot N, Philippe H. 2006. Computing Bayes factors using thermodynamic integration. *Systematic Biology*. 55:195–207.
- Lavine M, Schervish MJ. 1999. Bayes factors: What they are and what they are not. *American Statistician*. 53:119–122.
- McGuire JA, Witt CC, Altshuler DL, Remsen JV. 2007. Phylogenetic systematics and biogeography of hummingbirds: Bayesian and maximum likelihood analyses of partitioned data and selection of an appropriate partitioning strategy. *Systematic Biology*. 56:837–856.
- Nylander JAA, Ronquist F, Huelsenbeck JP, Aldrey JLN. 2004. Bayesian phylogenetic analysis of combined data. *Systematic Biology*. 53:47–67.
- Rannala B, Yang Z. 1996. Probability distribution of molecular evolutionary trees: A new method of phylogenetic inference. *Journal of Molecular Evolution*. 43:304–311.
- Robert CP, Casella G. 2002. *Monte Carlo Statistical Methods*. New York: Springer.
- Rodrigue N, Philippe H, Lartillot N. 2008. Uniformization for sampling realizations of Markov processes: Applications to Bayesian implementations of codon substitution models. *Bioinformatics*. 24:56–62.

- Ronquist F, Huelsenbeck JP. 2003. MrBayes 3: Bayesian phylogenetic inference under mixed models. *Bioinformatics*. 19:1572–1574.
- Ronquist F, van der Mark P, Huelsenbeck JP. 2009. Bayesian analysis of molecular evolution using MrBayes. In: Lemey P, Salemi M, Vandamme AM, editors, *The Phylogenetic Handbook: A Practical Approach to Phylogenetic Analysis and Hypothesis Testing*, Second Edition. Cambridge University Press, pp. 1–1.
- Rubinstein R. 1981. *Simulation and the Monte Carlo Method*. John Wiley & Sons, Inc. New York, NY, USA.
- Simon D, Larget B. 2001. Bayesian analysis in molecular biology and evolution (BAMBE). <http://www.mathcs.duq.edu/larget/bambe.html>.
- Smith A, Roberts G. 1993. Bayesian computation via the Gibbs sampler and related Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society. Series B (Methodological)*. 55:3–23.
- Suchard M, Weiss R, Sinsheimer J. 2005. Models for estimating bayes factors with applications to phylogeny and tests of monophyly. *Biometrics*. 61:665–673.
- Suchard MA, Weiss RE, Sinsheimer JS. 2001. Bayesian selection of continuous-time Markov chain evolutionary models. *Molecular Biology and Evolution*. 18:1001–1013.
- Sukumaran J, Holder MT. 2010. DendroPy: A Python library for phylogenetic computing. *Bioinformatics*. 26:1569–1571.
- Verdinelli I, Wasserman L. 1995. Computing Bayes factors using a generalization of the Savage Dickey density ratio. *Journal of the American Statistical Association*. 90:614–618.
- Wong KM, Suchard MA, Huelsenbeck JP. 2008. Alignment uncertainty and genomic analysis. *Science*. 319:473–476.
- Wu CH, Suchard MA, Drummond AJ. 2013. Bayesian selection of nucleotide substitution models and their site assignments. *Molecular Biology and Evolution*. 30:669–688.
- Xie W, Lewis PO, Fan Y, Kuo L, Chen MH. 2011. Improving marginal likelihood estimation for Bayesian phylogenetic model selection. *Systematic Biology*. 60:150–160.
- Yang Z. 2007. Fair-balance paradox, star-tree paradox, and Bayesian phylogenetics. *Molecular biology and evolution*. 24:1639.
- Yang Z, Rannala B. 1997. Bayesian phylogenetic inference using DNA sequences: A Markov chain Monte Carlo method. *Molecular Biology and Evolution*. 14:717–724.
- Yang Z, Rannala B. 2005. Branch-length prior influences Bayesian posterior probability of phylogeny. *Systematic Biology*. 54:455–470.