

UC Irvine

UC Irvine Previously Published Works

Title

Marginal Likelihoods in Phylogenetics: A Review of Methods and Applications

Permalink

<https://escholarship.org/uc/item/5q41f7gx>

Journal

Systematic Biology, 68(5)

ISSN

1063-5157

Authors

Oaks, Jamie R
Cobb, Kerry A
Minin, Vladimir N
[et al.](#)

Publication Date

2019-09-01

DOI

10.1093/sysbio/syz003

Copyright Information

This work is made available under the terms of a Creative Commons Attribution License, available at <https://creativecommons.org/licenses/by/4.0/>

Peer reviewed

Marginal likelihoods in phylogenetics: a review of methods and applications

Jamie R. Oaks^{*1}, Kerry A. Cobb¹, Vladimir N. Minin², and Adam D. Leaché³

¹Department of Biological Sciences & Museum of Natural History, Auburn University, Auburn, Alabama 36849

²Department of Statistics, University of California, Irvine, California 92697

³Department of Biology & Burke Museum of Natural History and Culture, University of Washington, Seattle, Washington 98195

May 11, 2018

Running head: Marginal likelihoods in phylogenetics

Abstract

By providing a framework of accounting for the shared ancestry inherent to all life, phylogenetics is becoming the statistical foundation of biology. The importance of model choice continues to grow as phylogenetic models continue to increase in complexity to better capture micro and macroevolutionary processes. In a Bayesian framework, the marginal likelihood is how data update our prior beliefs about models, which gives us an intuitive measure of comparing model fit that is grounded in probability theory. Given the rapid increase in the number and complexity of phylogenetic models, methods for approximating marginal likelihoods are increasingly important. Here we try to provide an intuitive description of marginal likelihoods and why they are important in Bayesian model testing. We also categorize and review methods for estimating marginal likelihoods of phylogenetic models. In doing so, we use simulations to evaluate the performance of one such method based on approximate-Bayesian computation (ABC) and find that it is biased as predicted by theory. Furthermore, we review some applications of marginal likelihoods to phylogenetics, highlighting how they can be used to learn about models of evolution from biological data. We conclude by discussing the challenges of Bayesian model choice and future directions that promise to improve the approximation of marginal likelihoods and Bayesian phylogenetics as a whole.

KEY WORDS: phylogenetics, marginal likelihood, model choice

*Corresponding author: joaks@auburn.edu

1 Introduction

Phylogenetics is rapidly progressing as the statistical foundation of comparative biology, providing a framework that accounts for the shared ancestry inherent in biological data. Soon after phylogenetics became feasible as a likelihood-based statistical endeavor (Felsenstein, 1981), models became richer to better capture processes of biological diversification and character change. This increasing trend in model complexity made Bayesian approaches appealing, because they can approximate posterior distributions of rich models by leveraging prior information and hierarchical models, where researchers can take into account uncertainty of all levels in the hierarchy.

From the earliest days of Bayesian phylogenetics (Rannala and Yang, 1996; Mau and Newton, 1997), the numerical tool of choice for approximating the posterior distribution was Markov chain Monte Carlo (MCMC). The popularity of MCMC was due, in no small part, to avoiding the calculation of the marginal likelihood of the model—the probability of the data under the model, averaged, with respect to the prior, over the whole parameter space. This marginalized measure of model fit is not easy to compute due to the large number of parameters in phylogenetic models over which the likelihood needs to be summed or integrated.

Nonetheless, marginal likelihoods are central to model comparison in a Bayesian framework. If we want to compare the fit of phylogenetic models, and in the process learn about evolutionary processes, we cannot avoid calculating marginal likelihoods. As the diversity and richness of phylogenetic models has increased, there has been a renewed appreciation of the importance of such Bayesian model comparison. As a result, there has been substantial work over the last decade to develop methods for estimating marginal likelihoods of phylogenetic models.

The goals of this review are to (1) try to provide some intuition about what marginal likelihoods are and why they are useful, (2) review the various methods available for approximating marginal likelihoods of phylogenetic models, (3) review some of the ways marginal likelihoods have been applied to learn about evolutionary history and processes, (4) discuss some of the challenges of Bayesian model choice, and (5) highlight some promising avenues for advancing the field of Bayesian phylogenetics.

2 What are marginal likelihoods and why are they useful?

A marginal likelihood is the average fit of a model to a dataset. More specifically, it is an average over the entire parameter space of the likelihood weighted by the prior. For a phylogenetic model M with parameters that include the discrete topology (T) and continuous branch lengths and other parameters that govern the evolution of the characters along the tree (together represented by θ), the marginal likelihood can be represented as

$$p(D | M) = \sum_T \int_{\theta} p(D | T, \theta, M) p(T, \theta | M) d\theta, \quad (1)$$

where D are the data. Each parameter of the model adds a dimension to the model, over which the likelihood must be averaged. The marginal likelihood is also the proportionality constant in the denominator of Bayes' rule that ensures the posterior is a proper probability density that sums and integrates to one:

$$p(T, \theta | D, M) = \frac{p(D | T, \theta, M)p(T, \theta | M)}{p(D | M)}. \quad (2)$$

Marginal likelihoods are the currency of model comparison in a Bayesian framework. The frequentist approach to model choice is based on comparing the maximum probability of the data under two models either using a likelihood ratio test or some information-theoretic criterion. Because adding a parameter (dimension) to a model will always ensure a maximum likelihood at least as large as without the parameter, some penalty must be imposed when parameters are added.

From a Bayesian perspective, we are interested in comparing the average fit of a model, rather than the maximum. This imposes a “natural” penalty for parameters, because each additional parameter introduces a dimension that must be averaged over. If that dimension introduces substantial parameter space with small likelihood, and little space that improves the likelihood, it will decrease the marginal likelihood. Thus, unlike the maximum likelihood, adding a parameter to a model can decrease the *marginal* likelihood.

The ratio of two marginal likelihoods gives us the factor by which the average fit of the model in the numerator is better or worse than the model in the denominator. This is called the Bayes factor (Jeffreys, 1935). We can again leverage Bayes' rule to gain more intuition for how marginal likelihoods and Bayes factors guide Bayesian model selection by writing it in terms of the posterior probability of a model, M_1 , among N candidate models:

$$p(M_1 | D) = \frac{p(D | M_1)p(M_1)}{\sum_{i=1}^N p(D | M_i)p(M_i)}. \quad (3)$$

This shows us that the posterior probability of a model is proportional to the prior probability multiplied by the marginal likelihood of that model. Thus, the marginal likelihood is how the data update our prior beliefs about a model. As a result, it is often simply referred to as “the evidence” (MacKay, 2005). If we look at the ratio of the posterior probabilities of two models,

$$\frac{p(M_1 | D)}{p(M_2 | D)} = \frac{p(D | M_1)}{p(D | M_2)} \times \frac{p(M_1)}{p(M_2)}, \quad (4)$$

we see that the Bayes factor is the factor by which the prior odds of a model is multiplied to give us the posterior odds. Thus, marginal likelihoods and their ratios give us intuitive measures of how much the data “favor” one model over another, and these measures have natural probabilistic interpretations.

2.1 A coin-flipping example

Before we discuss methods for approximating marginal likelihoods, let's use a simple, albeit contrived, example to help gain some intuition for marginal likelihoods and how they

differ from the posterior distribution of a model. Let’s assume we are interested in the probability of a coin we have not seen landing heads-side up when it is flipped (θ); we refer to this as the rate of landing heads up to avoid confusion with other uses of the word probability. Our plan is to flip this coin 100 times and count the number of times it lands heads up, which we model as a random outcome from a binomial distribution. Before flipping, we decide to compare four models that vary in our prior assumptions about the probability of the coin landing heads up (Figure 1): We assume

1. all values are equally probable ($M_1: \theta \sim \text{Beta}(1, 1)$),
2. the coin is likely weighted to land mostly “heads” or “tails” ($M_2: \theta \sim \text{Beta}(0.6, 0.6)$),
3. the coin is probably fair ($M_3: \theta \sim \text{Beta}(5.0, 5.0)$), and
4. the coin is weighted to land tails side up most of time ($M_4: \theta \sim \text{Beta}(1.0, 5.0)$).

We use beta distributions to represent our prior expectations, because the beta is a conjugate prior for the binomial likelihood function. This allows us to obtain the posterior distribution and marginal likelihood analytically.

After flipping the coin and observing that it landed heads side up 50 times, we can calculate the posterior probability distribution for the rate of landing heads up under each of our four models:

$$p(\theta | D, M_i) = \frac{p(D | \theta, M_i)p(\theta | M_i)}{p(D | M_i)}. \quad (5)$$

Doing so, we see that regardless of our prior assumptions about the rate of the coin landing heads, the posterior distribution is very similar (Figure 1). This makes sense; given we observed 50 heads out of 100 flips, values for θ toward zero and one are extremely unlikely, and the posterior is dominated by the likelihood of values near 0.5.

Given the posterior distribution for θ is very robust to our prior assumptions, we might assume that each of our four models explain the data similarly well. However, to compare their ability to explain the data, we need to average (integrate) the likelihood density function over all possible values of θ , weighting by the prior:

$$p(D | M_i) = \int_{\theta} p(D | \theta, M_i)p(\theta | M_i)d\theta. \quad (6)$$

Looking at the plots in Figure 1 we see that the models that place a lot of prior weight on values of θ that do not explain the data well (i.e., have small likelihood) have a much smaller marginal likelihood. Thus, even if we have very informative data that make the posterior distribution robust to prior assumptions, this example illustrates that the marginal likelihood of a model can still be very sensitive to the prior assumptions we make about the parameters. We have developed an interactive version of Figure 1 where readers can vary the parameters of the coin-flip experiment and prior assumptions to further gain intuition for marginal likelihoods (<https://kerrycobb.github.io/beta-binomial-web-demo/>). Now we turn to methods for approximating the marginal likelihood of phylogenetic models, where simple analytical solutions are not possible. Nonetheless, the same fundamental principles apply.

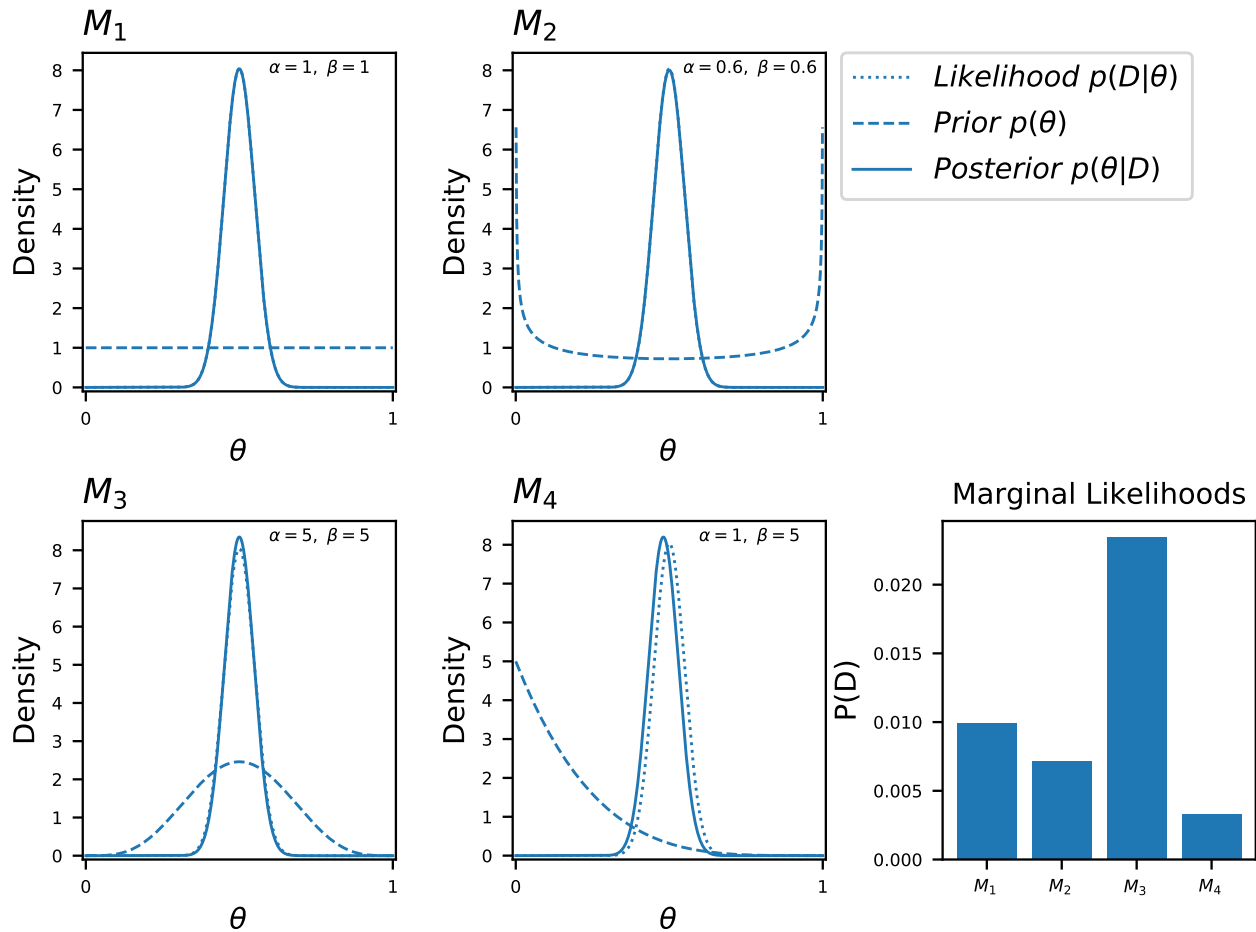


Figure 1. An illustration of the posterior probability densities and marginal likelihoods of the four different prior assumptions we made in our coin-flipping experiment. The data are 50 “heads” out of 100 coin flips, and the parameter, θ , is the probability of the coin landing heads side up. The binomial likelihood density function is proportional to a $\text{Beta}(51, 51)$ and is the same across the four different beta priors on θ (M_1 – M_4). The posterior of each model is a $\text{Beta}(\alpha + 50, \beta + 50)$ distribution. The marginal likelihoods ($P(D)$; the average of the likelihood density curve weighted by the prior) of the four models are compared.

3 Methods for marginal likelihood approximation

For all but the simplest of models, the summation and integrals in Equation 1 are analytically intractable. This is particularly true for phylogenetic models, which have a complex structure containing both discrete and continuous elements (Kim, 2000). Thus, we must resort to numerical techniques to approximate the marginal likelihood.

Perhaps the simplest numerical approximation of the marginal likelihood is to draw samples of a model’s parameters from their respective prior distributions. This turns the intractable integral into a sum of the samples’ likelihoods. Because the prior weight of each sample is one in this case, the marginal likelihood can be approximated by simply calculating the average likelihood of the prior samples. Alternatively, if we have a sample of the parameters from the posterior distribution—like one obtained from a “standard” Bayesian phylogenetic analysis via MCMC—we can again use summation to approximate the integral. In this case, the weight of each sample is the ratio of the prior density to the posterior density. As a result, the sum simplifies to the harmonic mean (HM) of the likelihoods from the posterior sample (Newton and Raftery, 1994). Both of these techniques can be thought of as importance-sampling integral approximations, and suffer from the fact that the prior and posterior are often *very* divergent, with the latter usually *much* more peaked than the former. A finite sample from the prior will often yield an underestimate of the marginal likelihood, because the region of parameter space with high likelihood is likely to be missed. Whereas a finite sample from the posterior will almost always lead to an overestimate (Lartillot and Philippe, 2006; Xie et al., 2011; Fan et al., 2011), because it will contain too few samples outside of the region of high likelihood, where the prior weight “penalizes” the average likelihood.

Recent methods developed to estimate marginal likelihoods generally fall into two categories for dealing with the sharp contrast between the prior and posterior that cripples the simple approaches mentioned above. One general strategy is to turn the giant leap between the unnormalized posterior and prior into many small steps. The second strategy is to turn the giant leap between the posterior and prior into a smaller leap between the posterior and a reference distribution that is as similar as possible to the posterior. These approaches are not mutually exclusive (e.g., see Fan et al. (2011)), but they serve as a useful way to categorize many of the methods available for approximating marginal likelihoods. In practical terms, the first strategy is computationally expensive, because samples need to be collected from each step between the posterior and prior, which is not normally part of a standard Bayesian phylogenetic analysis. The second strategy is very inexpensive because it attempts to approximate the marginal likelihood using only the posterior samples collected from a typical analysis.

3.1 Approaches that bridge the prior and posterior with small steps

3.1.1 Path sampling (PS)

Lartillot and Philippe (2006) introduced path sampling (also called thermodynamic integration) to phylogenetics to address the problem that the posterior is often dominated by the likelihood and very divergent from the prior. Rather than restrict themselves to a sample

from the posterior, they collected MCMC samples from a series of distributions between the prior and posterior. Specifically, samples are taken from a series of power-posterior distributions, where the likelihood is raised to a power β : $p(D|T, \theta, M)^\beta p(T, \theta | M)$. When $\beta = 1$, this is equal to the unnormalized joint posterior, which integrates to what we want to know, the marginal likelihood. When $\beta = 0$, this is equal to the joint prior distribution, which, assuming we are using proper prior probability distributions, integrates to 1. If we integrate the posterior expectation of the derivative with respect to β of the log power posterior over the interval (0–1) with respect to β , we get the ratio of the normalizing constants when β equals 1 and 0, and since we know the constant is 1 when β is zero, we are left with the marginal likelihood. [Lartillot and Philippe \(2006\)](#) approximated this integral by summing over MCMC samples taken from a discrete number of β values evenly distributed between 1 and 0.

3.1.2 Stepping stone (SS)

The stepping-stone method introduced by [Xie et al. \(2011\)](#) is similar to PS in that it also uses samples from power posteriors, but the motivation is not based on approximating the integral per se, but by the fact that we can accurately use importance sampling to approximate the ratio of normalizing constants at each step between the posterior and prior. Also, [Xie et al. \(2011\)](#) chose the values of β for the series of power posteriors from which to sample so that most were close to the prior (reference) distribution, rather than evenly distributed between 0 and 1. This is beneficial, because most of the change happens near the prior; the likelihood begins to dominate quickly, even at small values of β . The stepping-stone method results in more accurate estimates of the marginal likelihood with fewer steps than PS ([Xie et al., 2011](#)).

3.1.3 Generalized stepping stone (GSS)

The most accurate estimator of marginal likelihoods available to date, the generalized stepping-stone (GSS) method, leverages both strategies; taking many small steps from a starting point (reference distribution) that is much closer to the posterior than the prior ([Fan et al., 2011](#)). [Fan et al. \(2011\)](#) improved upon the original stepping-stone method by using a reference distribution that, in most cases, will be much more similar to the posterior than the prior. The reference distribution has the same form as the joint prior, but each marginal prior distribution is adjusted so that its mean and variance matches the corresponding sample mean and variance of an MCMC sample from the posterior. This guarantees that the support of the reference distribution will cover the posterior. Initially, the application of the GSS method was limited, because it required that the topology be fixed, because there was no reference distribution across topologies. However, [Holder et al. \(2014\)](#) introduced such a distribution on trees, allowing the GSS to approximate the fully marginalized likelihood of phylogenetic models.

3.1.4 Sequential Monte Carlo (SMC)

Another approach that uses sequential importance-sampling steps is sequential Monte Carlo (SMC), also known as particle filtering. Recently, SMC algorithms have been de-

veloped for approximating the posterior distribution of phylogenetic trees (Bouchard-Côté et al., 2012; Bouchard-Côté, 2014). While inferring the posterior, SMC algorithms can approximate the marginal likelihood of the model “for free,” by keeping a running average of the importance-sampling weights of the trees (particles) along the way. SMC algorithms hold a lot of promise for complementing MCMC in Bayesian phylogenetics due to their sequential nature and ease with which the computations can be parallelized (Bouchard-Côté et al., 2012; Dinh et al., 2016; Fourment et al., 2017). However, these approaches are still in their infancy, and the accuracy of SMC estimates of marginal likelihoods still need to be compared to the methods discussed above. See Bouchard-Côté (2014) for an accessible treatment of SMC in phylogenetics.

3.1.5 Nested sampling (NS)

Recently, Maturana R. et al. (2017) introduced the numerical technique known as nested sampling to Bayesian phylogenetics. This tries to simplify the multi-dimensional integral in Equation 1 into a one-dimensional integral over the cumulative distribution function of the likelihood. The latter can be numerically approximated using basic quadrature methods, essentially summing up the area of polygons under the likelihood function. The algorithm works by starting with a random sample of parameter values from the joint prior distribution and their associated likelihood scores. Sequentially, the sample with the lowest likelihood is removed and replaced by another random sample from the prior with the constraint that its likelihood must be larger than the removed sample. The approximate marginal likelihood is a running sum of the likelihood of these removed samples with appropriate weights. Re-sampling these removed samples according to their weights yields a posterior sample at no extra computational cost. Initial assessment of NS suggest it performs similarly to GSS. As with SMC, NS seems like a promising complement to MCMC for both approximating the posterior and marginal likelihood of phylogenetic models.

3.2 Approaches that use only posterior samples

3.2.1 Generalized harmonic mean (GHM)

Gelfand and Dey (1994) introduced a generalized harmonic mean estimator that uses an arbitrary normalized reference distribution, as opposed to the prior distribution used in the HM estimator, to weight the samples from the posterior. If the chosen reference distribution is more similar to the posterior than the prior (i.e., a “smaller leap” as discussed above), the GHM estimator will perform better than the HM estimator. However, for high-dimensional phylogenetic models, choosing a suitable reference distribution is very challenging, especially for tree topologies. As a result, the GHM estimator has not been used for comparing phylogenetic models. However, recent advances on defining a reference distribution on trees (Holder et al., 2014) makes the GHM a tenable option in phylogenetics.

3.2.2 Inflated-density ratio (IDR)

The inflated-density ratio estimator solves the problem of choosing a reference distribution by using a perturbation of the posterior density; essentially the posterior is “inflated”

from the center by a known radius (Petris and Tardella, 2007; Arima and Tardella, 2012, 2014). As one might expect, the radius must be chosen carefully. The application of this method to phylogenetics has been limited by the fact that all parameters must be unbounded; any parameters that are bounded (e.g., must be positive) must be re-parameterized to span the real number line. As a result, this method cannot be applied directly to MCMC samples collected by popular Bayesian phylogenetic software packages. Nonetheless, the IDR estimator has recently been applied to phylogenetic models (Arima and Tardella, 2014), including in settings where the topology is allowed to vary (Wu et al., 2014). Initial applications of the IDR are very promising, demonstrating comparable accuracy to methods that sample from power-posterior distributions while avoiding such computation (Arima and Tardella, 2014; Wu et al., 2014). Currently, however, the IDR has only been used on relatively small datasets and simple models of character evolution. More work is necessary to determine whether the promising combination of accuracy and computational efficiency holds for large datasets and rich models.

3.2.3 Partition-weighted kernel (PWK)

Recently, Wang et al. (2017) introduced the partition weighted kernel (PWK) method of approximating marginal likelihoods. This approach entails partitioning parameter space into regions within which the posterior density is relatively homogeneous. Given the complex structure of phylogenetic models, it is not obvious how this would be done. As of yet, this method has not been used for phylogenetic models. However, for simulations of mixtures of bivariate normal distributions, the PWK outperforms the IDR estimator (Wang et al., 2017). Thus, the method holds promise if it can be adapted to phylogenetic models.

3.3 Bayesian model averaging

An alternative to using marginal likelihoods for Bayesian model comparison is to sample across the competing models directly. The frequency of samples from each model approximates their posterior probability, which can be used to approximate Bayes factors among models. Algorithms for sampling across models include reversible-jump MCMC (Green, 1995), Gibbs sampling (Neal, 2000), Bayesian stochastic search variable selection (George and McCulloch, 1993; Kuo and Mallick, 1998), and approximations of reversible-jump (Jones et al., 2015). In fact, the first application of Bayes factors for phylogenetic model comparison was performed by Suchard et al. (2001) via reversible-jump MCMC. This technique was also used in Bayesian tests of phylogenetic incongruence/recombination (Suchard et al., 2003; Minin et al., 2005). In terms of selecting the correct “relaxed-clock” model from simulated data, Baele and Lemey (2014) showed that model-averaging performed similarly to the path-sampling and stepping-stone marginal likelihood estimators.

There are a couple of limitations for these approaches. First, a Bayes factor that includes a model with small posterior probability will suffer from Monte Carlo error. For example, unless a very large sample from the posterior is collected, some models might not be sampled at all. Second, and perhaps more importantly, for these numerical algorithms to be able to “jump” among models, the models being sampled need to be similar.

If two highly dissimilar models need to be compared, [Lartillot and Philippe \(2006\)](#) introduced a method of using path sampling to directly approximate the Bayes factor. Similarly, [Baele et al. \(2013\)](#) extended the stepping-stone approach of [Xie et al. \(2011\)](#) to do the same. However, if there are many models to compare, doing MCMC over power posteriors for every pairwise comparison will quickly become computationally prohibitive; approximating the marginal likelihood of each model would be simpler.

3.4 Approximate-likelihood approaches

Approximate-likelihood Bayesian computation (ABC) approaches ([Tavaré et al., 1997](#); [Beaumont et al., 2002](#)) have become popular in situations where it is not possible (or undesirable) to derive and compute the likelihood function of a model. The basic idea is simple: by generating simulations under the model, the fraction of times that we generate a simulated dataset that matches the observed data is a Monte Carlo approximation of the likelihood. Because simulating the observed data exactly is often not possible (or extremely unlikely), simulations “close enough” to the observed data are counted, and usually a set of insufficient summary statistics are used in place of the data. Whether a simulated dataset is “close enough” to count is formalized as whether or not it falls within a zone of tolerance around the empirical data.

This simple approach assumes the likelihood within the zone of tolerance is uniform. However, this zone usually needs to be quite large for computational tractability, so this assumption does not hold. [Leuenberger and Wegmann \(2010\)](#) proposed fitting a general linear model (GLM) to approximate the likelihood within the zone of tolerance. With the GLM in hand, the marginal likelihood of the model can simply be approximated by the marginal density of the GLM.

The accuracy of this estimator has not been assessed. However, there are good theoretical reasons to be skeptical of its accuracy. Because the GLM is only fit within the zone of tolerance (also called the “truncated prior”), it cannot account for the weight of the prior on the marginal likelihood outside of this region. Whereas the posterior distribution usually is not strongly influenced by regions of parameter space with low likelihood, the marginal likelihood very much is. By not accounting for prior weight in regions of parameter space outside the zone of tolerance, where the likelihood is low, we predict this method will not properly penalize models and tend to favor models with more parameters. This is analogous with how the harmonic mean estimator tends to overestimate the marginal likelihood due to having very few samples outside of the region of high likelihood.

To test this prediction, we assessed the behavior of the ABC-GLM method on 100 datasets simulated under the simplest possible phylogenetic model: two DNA sequences separated by a single branch along which the sequence evolved under a Jukes-Cantor model of nucleotide substitution ([Jukes and Cantor, 1969](#)). The simulated sequences were 10,000 nucleotides long, and the prior on the only parameter in the model, the length of the branch, was a uniform distribution from 0.0001 to 0.1 substitutions per site. For such a simple model, we used quadrature integration to calculate the marginal likelihood for each simulated alignment of two sequences. Integration using 1,000 and 10,000 steps and rectangular and trapezoidal quadrature rules all yielded identical values for the log marginal likelihood to at least five decimal places for all 100 simulated data sets, providing a very precise proxy for the true

values. We used a sufficient summary statistic, the proportion of variable sites, for ABC analyses. However, the ABC-GLM and quadrature marginal likelihoods are not directly comparable, because the marginal probability of the proportion of variable sites versus the site pattern counts will be on different scales that are data set dependent. So, we compare the ratio of marginal likelihoods (i.e., Bayes factors) comparing the correct branch-length model [branch length \sim uniform(0.0001, 0.1)] to a model with a prior approximately twice as broad [branch length \sim uniform(0.0001, 0.2)].

This very simple model is a good test of the ABC-GLM marginal likelihood estimator for several reasons. The use of a sufficient statistic for a finite, one-dimensional model makes ABC nearly equivalent to a full-likelihood Bayesian method (Figure S1). Thus, this is a “best-case scenario” for the ABC-GLM approach. Also, we can use quadrature integration for very good proxies for the true Bayes factors. Lastly, the simple scenario gives us some analytical expectations for the behavior of ABC-GLM. If it cannot penalize the marginal likelihood for the additional branch length space in the model with the broader prior, the Bayes factor should be off by a factor of approximately 2, or more precisely $(0.2 - 0.0001)/(0.1 - 0.0001)$. As shown in Figure 2, this is exactly what we find. This confirms our prediction that the ABC-GLM approach cannot average over regions of parameter space with low likelihood and thus will be biased toward favoring models with more parameter space. Given that the GLM approximation of the likelihood is only fit within a subset of parameter space with high likelihood, which is usually a *very* small region of a model, the marginal of the GLM should not be considered a marginal likelihood of the model. We want to emphasize that our findings in no way detract from the usefulness of ABC-GLM for parameter estimation.

Full details of these analyses, which were all designed atop the DendroPy phylogenetic API (version 4.3.0 commit 72ce015) (Sukumaran and Holder, 2010), can be found in the supplementary materials, and all of the code to replicate our results is freely available at <https://github.com/phyletica/abc-glm-marginal-test>.

4 Uses of marginal likelihoods

The application of marginal likelihoods to compare phylogenetic models is rapidly gaining popularity. Rather than attempt to be comprehensive, below we highlight examples that represent some of the diversity of questions being asked and the insights that marginal likelihoods can provide about our data and the evolutionary processes giving rise to them.

4.1 Comparing partitioning schemes

One of the earliest applications of marginal likelihoods in phylogenetics was to choose among ways of assigning models of substitution to different subsets of aligned sites. This became important when phylogenetics moved beyond single-locus trees to concatenated alignments of several loci. Mueller et al. (2004), Nylander et al. (2004), and Brandley et al. (2005) used Bayes factors calculated from harmonic mean estimates of marginal likelihoods to choose among different strategies for partitioning aligned characters to substitution models. All three studies found that the model with the most subsets was strongly preferred. Nylander et al. (2004) also showed that removing parameters for which the data seemed

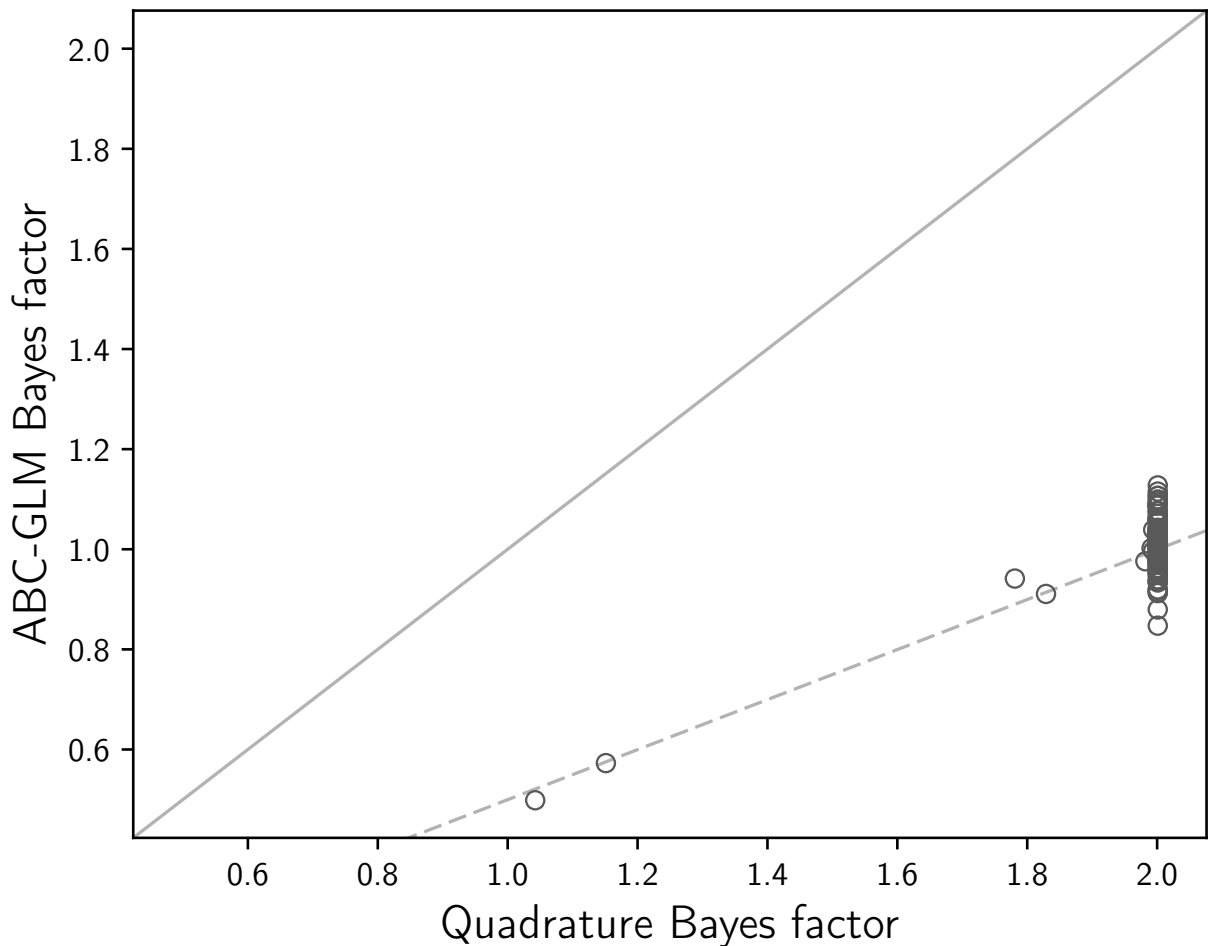


Figure 2. A comparison of the approximate-likelihood Bayesian computation general linear model (ABC-GLM) estimator of the marginal likelihood (Leuenberger and Wegmann, 2010) to quadrature integration approximations (Xie et al., 2011) for 100 simulated datasets. We compared the ratio of the marginal likelihood (Bayes factor) comparing the correct branch-length model [branch length $\sim \text{uniform}(0.0001, 0.1)$] to a model with a boader prior on the branch length [branch length $\sim \text{uniform}(0.0001, 0.2)$]. The solid line represents perfect performance of the ABC-GLM estimator (i.e., matching the “true” value of the Bayes factor). The dashed line represents the expected Bayes factor when failing to penalize for the extra parameter space (branch length 0.1 to 0.2) with essentially zero likelihood. Quadrature integration with 1,000 and 10,000 steps using the rectangular and trapezoidal rule produced identical values of log marginal likelihoods to at least five decimal places for all 100 simulated datasets.

to have little influence decreased the HM estimates of the marginal likelihood, suggesting that the HM estimates might favor over-parameterized models. These findings could be an artefact of the tendency of the HM estimator to overestimate marginal likelihoods and thus underestimate the “penalty” associated with the prior weight of additional parameters. However, [Brown and Lemmon \(2007\)](#) showed that for simulated data, HM estimates of Bayes factors can have a low error rate of over-partitioning an alignment.

[Fan et al. \(2011\)](#) showed that, again, the HM estimator strongly favors the most partitioned model for a four-gene alignment from cicadas (12 subsets partitioned by gene and codon position). However, the marginal likelihoods estimated via the generalized stepping stone method favor a much simpler model (3 subsets partitioned by codon position). This demonstrates the bias of simple importance-sampling methods when applied to finite samples from the posterior to estimate the average likelihood of phylogenetic models. It also suggests that relatively few, well-assigned subsets can go a long way to explain the variation in substitution rates among sites.

4.2 Comparing models of character substitution

[Lartillot and Philippe \(2006\)](#) used path sampling to compare models of amino-acid substitution. They found that the harmonic mean estimator favored the most parameter rich model for all five datasets they explored, whereas the path-sampling estimates favored simpler models for three of the datasets. This again demonstrates that accurately estimated marginal likelihoods can indeed “penalize” for over-parameterization of phylogenetic models. More importantly, this work also revealed that modeling heterogeneity in amino acid composition across sites of an alignment better explains the variation in biological data.

4.3 Comparing “relaxed clock” models

[Lepage et al. \(2007\)](#) used path sampling to approximate Bayes factors comparing various “relaxed-clock” phylogenetic models for three empirical datasets. They found that models in which the rate of substitution evolves across the tree (autocorrelated rate models) better explain the empirical sequence alignments they investigated than models that assume the rate of substitution on each branch is independent (uncorrelated rate models). This provides insight into how the rate of evolution evolves through time.

[Baele et al. \(2012\)](#) demonstrated that modeling among-branch rate variation with a lognormal distribution tends to explain mammalian sequence alignments better than using an exponential distribution. They used marginal likelihoods (PS and SS estimates) and Bayesian model averaging to compare the fit of lognormally and exponentially distributed relaxed clocks for almost 1,000 loci from 12 mammalian species. They found that the lognormal relaxed-clock was a better fit for almost 88% of the loci. [Baele et al. \(2012\)](#) also used marginal likelihoods to demonstrate the importance of using sampling dates when estimating time-calibrated phylogenetic trees. They used path-sampling and stepping-stone methods to estimate the marginal likelihoods of strict and relaxed-clock models for sequence data of herpes viruses. They found that when the dates the viruses were sampled were provided, a strict molecular clock was the best fit model, but when the dates were excluded, relaxed-clock models were strongly favored. Their findings show that using information about the

ages of the tips can be critical for accurately modeling processes of evolution and inferring evolutionary history.

4.4 Comparing demographic models

Baele et al. (2012) used the path-sampling and stepping-stone estimators for marginal likelihoods to compare the fit of various demographic models to the HIV-1 group M data of Worobey et al. (2008), and Methicillin-resistant *Staphylococcus aureus* (MRSA) data of Gray et al. (2011). They found that a flexible, nonparametric model that enforces no particular demographic history is a better explanation of the HIV and MRSA sequence data than exponential and logistic population growth models. This suggests that traditional parametric growth models are not the best predictors of viral and bacterial epidemics.

4.5 Measuring phylogenetic information content across genomic data sets

Not only can we use marginal likelihoods to learn about evolutionary models, but we can also use them to learn important lessons about our data. Brown and Thomson (2017) explored six different genomic data sets that were collected to infer phylogenetic relationships within Amniota. For each locus across all six data sets, they used the stepping-stone method (Xie et al., 2011) to approximate the marginal likelihood of models that included or excluded a particular branch (bipartition) in the amniote tree. This allowed Brown and Thomson (2017) to calculate, for each gene, Bayes factors as measures of support for or against particular relationships, some of which are uncontroversial (e.g., the monophyly of birds) and others contentious (e.g., the placement of turtles).

Their use of marginal likelihoods allowed Brown and Thomson (2017) to reveal a large degree of variation among loci in support for and against relationships that was masked by the corresponding posterior probabilities estimated by MCMC. Furthermore, they found that a small number of loci can have a large affect on the tree and associated posterior probabilities of branches inferred from the combined data. For example, they showed that including or excluding just two loci out of the 248 locus dataset of (Chiari et al., 2012) resulted in a posterior probability of 1.0 in support of turtles either being sister to crocodylians or archosaurs (crocodylians and birds), respectively. By using marginal likelihoods of different topologies, Brown and Thomson (2017) were able to identify these two loci as putative paralogs due to their strikingly strong support for turtles being sister to crocodylians. This work demonstrates how marginal likelihoods can simultaneously be used as a powerful means of controlling the quality of data in “phylogenomics”, while informing us about the evolutionary processes that gave rise to our data.

Furthermore, Brown and Thomson (2017) found that the properties of loci commonly used as proxies for the reliability of phylogenetic signal (rate of substitution, how “clock-like” the rate is, base composition heterogeneity, amount of missing data, and alignment uncertainty) were poor predictors of Bayes factor support for well-established amniote relationships. This suggests these popular rules of thumb are not useful for identifying “good” loci for phylogenetic inference.

4.6 Phylogenetic factor analysis

The goal of comparative biology is to understand the relationships among a potentially large number of phenotypic traits across organisms. To do so correctly, we need to account for the inherent shared ancestry underlying all life (Felsenstein, 1985). A lot of progress has been made for inferring the relationship between pairs of phenotypic traits as they evolve across a phylogeny, but a general and efficient solution for large numbers of continuous and discrete traits has remained elusive. Tolkoﬀ et al. (2017) introduced Bayesian factor analysis to a phylogenetic framework as a potential solution. Phylogenetic factor analysis works by modeling a small number of unobserved (latent) factors that evolve independently across the tree, which give rise to the large number of observed continuous and discrete phenotypic traits. This allows correlations among traits to be estimated, without having to model every trait as a conditionally independent process.

The question that immediately arises is, what number of factors best explains the evolution of the observed traits? To address this, (Tolkoﬀ et al., 2017) use path sampling to approximate the marginal likelihood of models with different numbers of traits. To do so, they extend the path sampling method to handle the latent variables underlying the discrete traits by softening the thresholds that delimit the discrete character states across the series of power posteriors. This new approach leverages Bayesian model comparison via marginal likelihoods to learn about the processes governing the evolution of multidimensional phenotypes.

4.7 Comparing phylogeographic models

Phylogeographers are interested in explaining the genetic variation within and among species across a landscape. As a result, we are often interested in comparing models that include various combinations of micro and macro-evolutionary processes and geographic and ecological parameters. Deriving the likelihood function for such models is often diﬃcult and, as a result, model choice approaches that use approximate-likelihood Bayesian computation (ABC) are often used.

At the forefront of generalizing phylogeographic models is an approach that is referred to as iDDC, which stands for integrating distributional, demographic, and coalescent models (Papadopoulou and Knowles, 2016). This approach simulates data under various phylogeographical models upon proxies for habitat suitability derived from species distribution models. To choose the model the best explains the empirical data, this approach uses the marginal densities of the models estimated via the ABC-GLM method and p-values derived from these densities (He et al., 2013) (Massatti and Knowles, 2016) (Bemmels et al., 2016) (Knowles and Massatti, 2017) (Papadopoulou and Knowles, 2016). This approach is an important step forward for bringing more biological realism into phylogeographical models. However, given that the marginal GLM density fitted to a truncated region of parameter space should not be interpreted as a marginal likelihood of the full model (see above; Figure 2), these methods should be seen as a useful exploration of data, rather than rigorous hypothesis tests. Because ABC-GLM marginal densities fail to penalize parameters for their prior weight in regions of low likelihood, these approaches will likely be biased toward over-parameterized phylogeographical models. Nonetheless, knowledge of this bias can help guide

interpretations of results.

4.8 Species delimitation

Calculating the marginal probability of sequence alignments (Grummer et al., 2013) and single-nucleotide polymorphisms (Leaché et al., 2014) under various multi-species coalescent models has been used to estimate species boundaries. By comparing the marginal likelihoods of models that differ in how they assign individual organisms to species, systematists can calculate Bayes factors to determine how much the genetic data support different delimitations. Using simulated data, (Grummer et al., 2013) found that marginal likelihoods calculated using path sampling and stepping-stone methods outperformed harmonic mean estimators at identifying the true species delimitation model. Marginal likelihoods seem better able to distinguish some species delimitation models than others. For example, models that lump species together or reassign samples into different species produce larger marginal likelihood differences versus models that split populations apart (Grummer et al., 2013; Leaché et al., 2014). Current implementations of the multi-species coalescent assume strict models of genetic isolation, and oversplitting populations that exchange genes creates a difficult Bayesian model comparison problem that does not include the correct model (Leaché et al. in prep.).

Species delimitation using marginal likelihoods in conjunction with Bayes factors has some advantages over alternative approaches. The flexibility of being able to compare non-nested models that contain different numbers of species, or different species assignments, is one key advantage. The methods also integrate over gene trees, species trees, and other model parameters, allowing the marginal likelihoods of delimitations to be compared without conditioning on any parameters being known. Marginal likelihoods also provide a natural way to rank competing models while automatically accounting for model complexity (Baele et al., 2012). Finally, it is unnecessary to assign prior probabilities to the alternative species delimitation models being compared. The marginal likelihood of a delimitation provides the factor by which the data update our prior expectations, regardless of what that expectation is (Equation 3). As multi-species coalescent models continue to advance, using the marginal likelihoods of delimitations will continue to be a powerful approach to learning about biodiversity.

5 Discussion

5.1 Promising future directions

As Bayesian phylogenetics continues to explore more complex models of evolution, and datasets continue to get larger, accurate and efficient methods of estimating marginal likelihoods will become increasingly important. Thanks to substantial work in recent years, robust methods have been developed, such as the generalized stepping-stone approach (Fan et al., 2011). However, these methods are computationally demanding as they have to sample likelihoods across a series of power-posterior distributions that are not useful for parameter estimation. Recent work has introduced promising methods to estimate marginal likelihoods solely from samples from the posterior distribution. However, these methods remain difficult

to apply to phylogenetic models, and their performance on rich models and large datasets remains to be explored.

Promising avenues for future research on methods for estimating marginal likelihoods of phylogenetic models include continued work on reference distributions that are as similar to the posterior as possible, but easy to formulate and use. This would improve the performance and applicability of the GSS and derivations of the GHM approach. Currently, the most promising method that works solely from a posterior sample is IDR. Making this method easier to apply to phylogenetic models and implementing it in popular Bayesian phylogenetic software packages, like RevBayes (Höhna et al., 2016) and BEAST (Drummond et al., 2012; Bouckaert et al., 2014) would be very useful, though nontrivial.

Furthermore, nested sampling and sequential Monte Carlo are exciting numerical approaches to Bayesian phylogenetics. These methods essentially use the same amount of computation to both sample from the posterior distribution of phylogenetic models and provide an approximation of the marginal likelihood. Both approaches are new to phylogenetics, but hold a lot of promise for Bayesian phylogenetics generally and model comparison via marginal likelihoods specifically.

5.2 The state of ABC approaches to Bayesian model choice

We found ABC estimation of marginal likelihoods can be problematic, even for the simplest of phylogenetic models (Figure 2). Other recent work has also demonstrated the poor performance of ABC methods for choosing among models (Robert et al., 2011; Oaks et al., 2013, 2014). Taken together, we recommend that Bayesian model choice based on approximate-likelihoods should be treated as a useful means of exploring data, rather than a robust statistical framework.

5.3 A fundamental challenge of Bayesian model choice

While the computational challenges to approximating marginal likelihoods are very real and will provide fertile ground for future research, it is often easy to forget about a fundamental challenge of Bayesian model choice. This challenge becomes apparent when we reflect on the differences between Bayesian model choice and parameter estimation. The posterior distribution of a model, and associated parameter estimates, are informed by the likelihood function (Equation 2), whereas the posterior probability of that model is informed by the *marginal* likelihood (Equation 3). When we have informative data, the posterior distribution is dominated by the likelihood, and as a result our parameter estimates are often robust to prior assumptions we make about the parameters. However, when comparing models, we need to assess their overall ability to predict the data, which entails averaging over the entire parameter space of the model, not just the regions of high likelihood. As a result, marginal likelihoods and associated model choices can be very sensitive to priors on the *parameters* of each model, even when the data are very informative (Figure 1). This sensitivity to prior assumptions about parameters is inherent to Bayesian model choice. Accordingly, the results of any application of Bayesian model selection should be accompanied by an assessment of the sensitivity of those results to the priors placed on the models' parameters.

5.4 Conclusions

Marginal likelihoods are intuitive measures of model fit that are grounded in probability theory. As a result, they provide us with a coherent way of gaining a better understanding about how evolution proceeds as we accrue biological data. We highlighted how marginal likelihoods of phylogenetic models can be used to learn about evolutionary processes and how our data inform our models. Because shared ancestry is a fundamental property of life, the use of marginal likelihoods of phylogenetic models promises to continue to advance biology.

6 Funding

This work was supported by the National Science Foundation (grant numbers DBI 1308885 and DEB 1656004 to JRO).

7 Acknowledgments

We thank Mark Holder for helpful discussions about comparing approximate and full marginal likelihoods. We also thank the members of the Phyletica Lab (the phyleticians) for helpful comments that improved an early draft of this paper. The computational work was made possible by the Auburn University (AU) Hopper Cluster supported by the AU Office of Information Technology. This paper is contribution number **XXXX** of the Auburn University Museum of Natural History.

References

- Arima, S. and L. Tardella. 2012. Improved harmonic mean estimator for phylogenetic model evidence. *Journal of Computational Biology* 19:418–438.
- Arima, S. and L. Tardella. 2014. Inflated density ratio (IDR) method for estimating marginal likelihoods in Bayesian phylogenetics. chap. 3, Pages 25–57 *in* *Bayesian phylogenetics: methods, algorithms, and applications* (M.-H. Chen, L. Kuo, and P. O. Lewis, eds.). CRC Press, Boca Raton, Florida, USA.
- Baele, G. and P. Lemey. 2014. Bayesian model selection in phylogenetics and genealogy-based population genetics. chap. 4, Pages 59–93 *in* *Bayesian phylogenetics: methods, algorithms, and applications* (M.-H. Chen, L. Kuo, and P. O. Lewis, eds.). CRC Press, Boca Raton, Florida, USA.
- Baele, G., P. Lemey, T. Bedford, A. Rambaut, M. A. Suchard, and A. V. Alekseyenko. 2012. Improving the accuracy of demographic and molecular clock model comparison while accommodating phylogenetic uncertainty. *Molecular Biology and Evolution* 29:2157–2167.

- Baele, G., P. Lemey, and S. Vansteelandt. 2013. Make the most out of your samples: Bayes factor estimators for high-dimensional models of sequence evolution. *BMC Bioinformatics* 14:85.
- Beaumont, M., W. Zhang, and D. J. Balding. 2002. Approximate Bayesian computation in population genetics. *Genetics* 162:2025–2035.
- Bemmels, J. B., P. O. Title, J. Ortego, and L. L. Knowles. 2016. Tests of species-specific models reveal the importance of drought in postglacial range shifts of a mediterranean-climate tree: insights from integrative distributional, demographic and coalescent modelling and ABC model selection. *Molecular Ecology* 25:4889–4906.
- Bouchard-Côté, A. 2014. SMC (sequential Monte Carlo) for bayesian phylogenetics. chap. 8, Pages 163–185 *in* Bayesian phylogenetics: methods, algorithms, and applications (M.-H. Chen, L. Kuo, and P. O. Lewis, eds.). CRC Press, Boca Raton, Florida, USA.
- Bouchard-Côté, A., S. Sankararaman, and M. I. Jordan. 2012. Phylogenetic inference via sequential Monte Carlo. *Systematic Biology* 61:579–93.
- Bouckaert, R., J. Heled, D. Kühnert, T. Vaughan, C.-H. Wu, D. Xie, M. A. Suchard, A. Rambaut, and A. J. Drummond. 2014. BEAST 2: A software platform for Bayesian evolutionary analysis. *PLOS Computational Biology* 10:1–6.
- Brandley, M. C., A. Schmitz, and T. W. Reeder. 2005. Partitioned Bayesian analyses, partition choice, and the phylogenetic relationships of scincid lizards. *Systematic Biology* 54:373–390.
- Brown, J. M. and A. R. Lemmon. 2007. The importance of data partitioning and the utility of Bayes factors in Bayesian phylogenetics. *Systematic Biology* 56:643–655.
- Brown, J. M. and R. C. Thomson. 2017. Bayes factors unmask highly variable information content, bias, and extreme influence in phylogenomic analyses. *Systematic Biology* 66:517–530.
- Chiari, Y., V. Cahais, N. Galtier, and F. Delsuc. 2012. Phylogenomic analyses support the position of turtles as the sister group of birds and crocodiles (Archosauria). *BMC Biology* 10:65.
- Dinh, V., A. E. Darling, and F. A. Matsen IV. 2016. Online Bayesian phylogenetic inference: theoretical foundations via Sequential Monte Carlo. [arXiv:1610.08148](https://arxiv.org/abs/1610.08148) [q-bio.PE] .
- Drummond, A. J., M. A. Suchard, D. Xie, and A. Rambaut. 2012. Bayesian phylogenetics with BEAUti and the BEAST 1.7. *Molecular Biology And Evolution* 29:1969–1973.
- Fan, Y., R. Wu, M.-H. Chen, L. Kuo, and P. O. Lewis. 2011. Choosing among partition models in Bayesian phylogenetics. *Molecular Biology And Evolution* 28:523–532.
- Felsenstein, J. 1981. Evolutionary trees from DNA sequences: A maximum likelihood approach. *Journal of Molecular Evolution* 17:368–376.

- Felsenstein, J. 1985. Phylogenies and the comparative method. *The American Naturalist* 125:1–15.
- Fourment, M., B. C. Claywell, V. Dinh, C. McCoy, F. A. Matsen IV, and A. E. Darling. 2017. Effective online Bayesian phylogenetics via sequential Monte Carlo with guided proposals. [bioRxiv](https://arxiv.org/abs/1708.02032) .
- Gelfand, A. E. and D. K. Dey. 1994. Bayesian model choice: asymptotic and exact calculations. *Journal of the Royal Statistical Society Series B* 56:501–514.
- George, E. I. and R. E. McCulloch. 1993. Variable selection via Gibbs sampling. *Journal of the American Statistical Association* 88:881–889.
- Gray, R. R., A. J. Tatem, J. A. Johnson, A. V. Alekseyenko, O. G. Pybus, M. A. Suchard, and M. Salemi. 2011. Testing spatiotemporal hypothesis of bacterial evolution using methicillin-resistant *Staphylococcus aureus* ST239 genome-wide data within a Bayesian framework. *Molecular Biology and Evolution* 28:1593–1603.
- Green, P. J. 1995. Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika* 82:711–732.
- Grummer, J. A., R. W. Bryson Jr., and T. W. Reeder. 2013. Species delimitation using Bayes factors: simulations and application to the *Sceloporus scalaris* species group (Squamata: Phrynosomatidae). *Systematic Biology* 63:119–133.
- He, Q., D. L. Edwards, and L. L. Knowles. 2013. Integrative testing of how environments from the past to the present shape genetic structure across landscapes. *Evolution* 67:3386–3402.
- Höhna, S., M. J. Landis, T. A. Heath, B. Boussau, N. Lartillot, B. R. Moore, J. P. Huelsenbeck, and F. Ronquist. 2016. RevBayes: Bayesian phylogenetic inference using graphical models and an interactive model-specification language. *Systematic Biology* 65:726–736.
- Holder, M. T., P. O. Lewis, D. L. Swofford, and D. Bryant. 2014. Variable tree topology stepping-stone marginal likelihood estimation. chap. 5, Pages 95–110 *in* Bayesian phylogenetics: methods, algorithms, and applications (M.-H. Chen, L. Kuo, and P. O. Lewis, eds.). CRC Press, Boca Raton, Florida, USA.
- Jeffreys, H. 1935. Some tests of significance, treated by the theory of probability. *Proceedings of the Cambridge Philosophical Society* 31:203–222.
- Jones, G., Z. Aydin, and B. Oxelman. 2015. DISSECT: An assignment-free Bayesian discovery method for species delimitation under the multispecies coalescent. *Bioinformatics* 31:991–998.
- Jukes, T. H. and C. R. Cantor. 1969. Evolution of protein molecules. chap. 24, Pages 21–132 *in* Mammalian Protein Metabolism (H. N. Munro, ed.) vol. III. Academic Press, New York.

- Kim, J. 2000. Slicing hyperdimensional oranges: the geometry of phylogenetic estimation. *Molecular Phylogenetics and Evolution* 17:58–75.
- Knowles, L. L. and R. Massatti. 2017. Distributional shifts—not geographic isolation—as a probable driver of montane species divergence. *Ecography* .
- Kuo, L. and B. Mallick. 1998. Variable selection for regression models. *Sankhyā: The Indian Journal of Statistics, Series B* 60:65–81.
- Lartillot, N. and H. Philippe. 2006. Computing Bayes factors using thermodynamic integration. *Systematic Biology* 55:195–207.
- Leaché, A. D., M. K. Fujita, V. N. Minin, and R. R. Bouckaert. 2014. Species delimitation using genome-wide SNP data. *Systematic Biology* 63:534–542.
- Lepage, T., D. Bryant, H. Philippe, and N. Lartillot. 2007. A general comparison of relaxed molecular clock models. *Molecular Biology And Evolution* 24:2669–2680.
- Leuenberger, C. and D. Wegmann. 2010. Bayesian computation and model selection without likelihoods. *Genetics* 184:243–252.
- MacKay, D. J. C. 2005. *Information Theory, Inference & Learning Algorithms*. 7.2 ed. Cambridge University Press, New York, New York, USA.
- Massatti, R. and L. L. Knowles. 2016. Contrasting support for alternative models of genomic variation based on microhabitat preference: species-specific effects of climate change in alpine sedges. *Molecular Ecology* 25:3974–3986.
- Maturana R., P., B. J. Brewer, and S. Klaere. 2017. Model selection and parameter inference in phylogenetics using nested sampling. [arXiv:1703.05471](https://arxiv.org/abs/1703.05471) [q-bio.QM] Pages 1–20.
- Mau, B. and M. A. Newton. 1997. Phylogenetic inference for binary data on dendrograms using Markov chain Monte Carlo. *Journal of Computational and Graphical Statistics* 6:122–131.
- Minin, V. N., K. S. Dorman, F. Fang, and M. A. Suchard. 2005. Dual multiple change-point model leads to more accurate recombination detection. *Bioinformatics* 21:3034–3042.
- Mueller, R. L., J. R. Macey, M. Jaekel, D. B. Wake, and J. L. Boore. 2004. Morphological homoplasy, life history evolution, and historical biogeography of plethodontid salamanders inferred from complete mitochondrial genomes. *Proceedings of the National Academy of Sciences of the United States of America* 101:13820–13825.
- Neal, R. M. 2000. Markov chain sampling methods for Dirichlet process mixture models. *Journal of Computational and Graphical Statistics* 9:249–265.
- Newton, M. A. and A. E. Raftery. 1994. Approximate Bayesian inference with the weighted likelihood bootstrap. *Journal of the Royal Statistical Society Series B-Statistical Methodology* 56:3–48.

- Nylander, J. A. A., F. Ronquist, J. P. Huelsenbeck, and J. L. Nieves-Aldrey. 2004. Bayesian phylogenetic analysis of combined data. *Systematic Biology* 53:47–67.
- Oaks, J. R., C. W. Linkem, and J. Sukumaran. 2014. Implications of uniformly distributed, empirically informed priors for phylogeographical model selection: A reply to hickerson et al. *Evolution* 68:3607–3617.
- Oaks, J. R., J. Sukumaran, J. A. Esselstyn, C. W. Linkem, C. D. Siler, M. T. Holder, and R. M. Brown. 2013. Evidence for climate-driven diversification? a caution for interpreting ABC inferences of simultaneous historical events. *Evolution* 67:991–1010.
- Papadopoulou, A. and L. L. Knowles. 2016. Toward a paradigm shift in comparative phylogeography driven by trait-based hypotheses. *Proceedings of the National Academy of Sciences* 113:8018–8024.
- Petris, G. and L. Tardella. 2007. New perspectives for estimating normalizing constants via posterior simulation. Tech. rep. Sapienza Università di Roma Roma, Italy.
- 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.
- Robert, C. P., J.-M. Cornuet, J.-M. Marin, and N. S. Pillai. 2011. Lack of confidence in approximate Bayesian computation model choice. *Proceedings of the National Academy of Sciences* 108:15112–15117.
- Suchard, M. A., R. E. Weiss, K. S. Dorman, and J. S. Sinsheimer. 2003. Inferring spatial phylogenetic variation along nucleotide sequences. *Journal of the American Statistical Association* 98:427–437.
- Suchard, M. A., R. E. Weiss, and J. S. Sinsheimer. 2001. Bayesian selection of continuous-time Markov chain evolutionary models. *Molecular Biology And Evolution* 18:1001–1013.
- Sukumaran, J. and M. T. Holder. 2010. DendroPy: a Python library for phylogenetic computing. *Bioinformatics* 26:1569–1571.
- Tavaré, S., D. J. Balding, R. C. Griffiths, and P. Donnelly. 1997. Inferring coalescence times from DNA sequence data. *Genetics* 145:505–518.
- Tolkoff, M. R., M. E. Alfaro, G. Baele, P. Lemey, and M. A. Suchard. 2017. Phylogenetic factor analysis. [arXiv:1701.07496](https://arxiv.org/abs/1701.07496) [stat.ME] .
- Wang, Y.-B., M.-H. Chen, L. Kuo, and P. O. Lewis. 2017. A new Monte Carlo method for estimating marginal likelihoods. *Bayesian Analysis* Pages 1–23.
- Wegmann, D., C. Leuenberger, S. Neuenchwander, and L. Excoffier. 2010. ABCtoolbox: a versatile toolkit for approximate bayesian computations. *BMC Bioinformatics* 11:116.

- Worobey, M., M. Gemmel, D. E. Teuwen, T. Haselkorn, K. Kunstman, M. Bunce, J.-J. Muyembe, J.-M. M. Kabongo, R. M. Kalengayi, E. Van Marck, M. T. P. Gilbert, and S. M. Wolinsky. 2008. Direct evidence of extensive diversity of HIV-1 in Kinshasa by 1960. *Nature* 455:661–664.
- Wu, R., M.-H. Chen, L. Kuo, and P. O. Lewis. 2014. Consistency of marginal likelihood estimation when topology varies. chap. 6, Pages 113–127 *in* Bayesian phylogenetics: methods, algorithms, and applications (M.-H. Chen, L. Kuo, and P. O. Lewis, eds.). CRC Press, Boca Raton, Florida, USA.
- Xie, W., P. O. Lewis, Y. Fan, L. Kuo, and M.-H. Chen. 2011. Improving marginal likelihood estimation for Bayesian phylogenetic model selection. *Systematic Biology* 60:150–160.

Figure Captions

- Figure 1. An illustration of the posterior probability densities and marginal likelihoods of the four different prior assumptions we made in our coin-flipping experiment. The data are 50 “heads” out of 100 coin flips, and the parameter, θ , is the probability of the coin landing heads side up. The binomial likelihood density function is proportional to a $\text{Beta}(51, 51)$ and is the same across the four different beta priors on θ (M_1 – M_4). The posterior of each model is a $\text{Beta}(\alpha + 50, \beta + 50)$ distribution. The marginal likelihoods ($P(D)$); the average of the likelihood density curve weighted by the prior) of the four models are compared.
- Figure 2. A comparison of the approximate-likelihood Bayesian computation general linear model (ABC-GLM) estimator of the marginal likelihood (Leuenberger and Wegmann, 2010) to quadrature integration approximations (Xie et al., 2011) for 100 simulated datasets. We compared the ratio of the marginal likelihood (Bayes factor) comparing the correct branch-length model [branch length $\sim \text{uniform}(0.0001, 0.1)$] to a model with a boader prior on the branch length [branch length $\sim \text{uniform}(0.0001, 0.2)$]. The solid line represents perfect performance of the ABC-GLM estimator (i.e., matching the “true” value of the Bayes factor). The dashed line represents the expected Bayes factor when failing to penalize for the extra parameter space (branch length 0.1 to 0.2) with essentially zero likelihood. Quadrature integration with 1,000 and 10,000 steps using the rectangular and trapezoidal rule produced identical values of log marginal likelihoods to at least five decimal places for all 100 simulated datasets.
- Figure S1. A comparison of the true branch length separating each pair of simulated sequences to the branch length estimated by ABC-GLM and full-likelihood MCMC under the correct branch-length model (branch length $\sim \text{uniform}(0.0001, 0.1)$) and the vague model (branch length $\sim \text{uniform}(0.0001, 0.1)$).

Supporting Information

Title: Marginal likelihoods in phylogenetics: a review of methods and applications

Authors: Jamie R. Oaks Corresponding author: joaks@auburn.edu¹, Kerry A. Cobb¹, Vladimir N. Minin², and Adam D. Leaché³

¹Department of Biological Sciences & Museum of Natural History, Auburn University, Auburn, Alabama 36849

²Department of Statistics, University of California, Irvine, California 92697

³Department of Biology & Burke Museum of Natural History and Culture, University of Washington, Seattle, Washington 98195

We set up a simple scenario for assessing the performance of the method for estimating marginal likelihoods based on approximating the likelihood function with a general linear model (GLM) fitted to posterior samples collected via approximate-likelihood Bayesian computation (ABC) (Leuenberger and Wegmann, 2010); hereforth referred to as ABC-GLM. The scenario is a DNA sequence, 10,000-nucleotides in length, that evolves along a branch according to a Jukes-Cantor continuous-time markov chain (CTMC) model of nucleotide substitution (Jukes and Cantor, 1969). Because the Jukes-Cantor model forces the relative rates of change among the four nucleotides and the equilibrium nucleotide frequencies to be equal, there is only a single parameter in the model, the length of the branch, and the direction of evolution along the branch does not matter.

1 Simulating data sets

We simulated 100 data sets under this model by

1. drawing 10,000 nucleotides of the “ancestral” sequence from their equilibrium frequencies ($\frac{1}{4}$),
2. drawing a branch length \sim uniform(0.0001, 0.1), and
3. evolving the sequence along the branch according to the Jukes-Cantor CTMC model to get the “descendant” sequence.

This was done using the DendroPy phylogenetic API (version 4.3.0 commit 72ce015) (Sukumar and Holder, 2010).

2 Calculating “true” Bayes factors

For each data set, we used quadrature approaches to approximate the marginal likelihood by integrating the posterior density over the branch length prior. We did this for two models:

1. the correct model [branch length \sim uniform(0.0001, 0.1)], and
2. a model with a branch length prior slightly more than twice as broad [branch length \sim uniform(0.0001, 0.2)], which we refer to as the “vague model”.

For both models and for each dataset we used the rectangular and trapezoidal quadrature rules with 1,000 and 10,000 steps (i.e., four approximations of the marginal likelihood for each data set under each model). Across all 100 data sets and both models, all four approximations were identical to at least five decimal places. For each data set, we calculated the log Bayes factor comparing the correct model to the vague model.

3 Approximate-likelihood Bayesian computation

To collect an approximate posterior sample from the correct model for a data set, we first calculated the proportion of variable sites ($Pvar$) between the two sequences. Next, we simulated 50,000 datasets under the correct model, calculated $Pvar$ for each of them, and retained the 1,000 samples with the values of $Pvar$ closest to that calculated from the data. Lastly, we used ABCtoolbox version 1.1 [Wegmann et al. \(2010\)](#) to fit a GLM to the retained samples and calculate the marginal density of the GLM, using a bandwidth of 0.002. We did the same to obtain an ABC-GLM estimate of the marginal density for the vague model with two differences: (1) we drew the branch length for each prior sample from the vague prior [branch length \sim uniform(0.0001, 0.2)], and (2) to maintain the same expected tolerance under both models, we simulated 100,000 datasets under the vague model (retaining the 1,000 samples closest to the $Pvar$ of the data).

For each data set, we calculated the log Bayes factor from the GLM marginal densities of the correct and vague model, and compared the ABC-GLM-estimated Bayes factor to the “true” Bayes factor calculated via quadrature integration (Figure 2).

4 Full-likelihood Markov chain Monte Carlo analyses

One goal of the simplicity of the above model is that the additional approximation of the ABC approach would be limited. All numerical Bayesian analyses, based on full or approximate likelihoods, suffer from Monte Carlo error associated with approximating the posterior with a finite number of samples. Approximate-likelihood methods usually suffer from two additional sources of approximation: (1) the full data are replaced with insufficient summary statistics, and (2) samples are retained that do not exactly match the data or summary statistics (i.e., the “tolerance” of ABC). In our analyses described above, we avoided the former source of error by using a sufficient statistic. We hoped to minimize the latter source of error by evaluating many samples from a one-dimensional model with finite bounds; we also kept this source of error approximately equal for both models by sampling in proportion to the width of the model.

To verify that the error introduced by the tolerance of the ABC analyses was minimal, we compared the branch length estimates to those estimated by full-likelihood Markov chain Monte Carlo (MCMC). For each data set, under both models, we ran a chain for 10,000 generations, sampling every 10 generations. All chains appeared to reach stationarity by the first sample (10th generation). We plotted the branch length estimated via ABC-GLM and MCMC under both the true and vague models against the true branch lengths. The results of all four analyses across all 100 data sets are almost indistinguishable (Figure S1), confirming that the approximation introduced by the tolerance is very minimal. Our ABC-GLM analyses are essentially equivalent to full-likelihood Bayesian analyses, creating a “best-case scenario” for evaluating the marginal likelihood estimates of the ABC-GLM method.

5 Reproducibility

All of the code to replicate our results is freely available at <https://github.com/phyletica/abc-glm-marginal-test>.

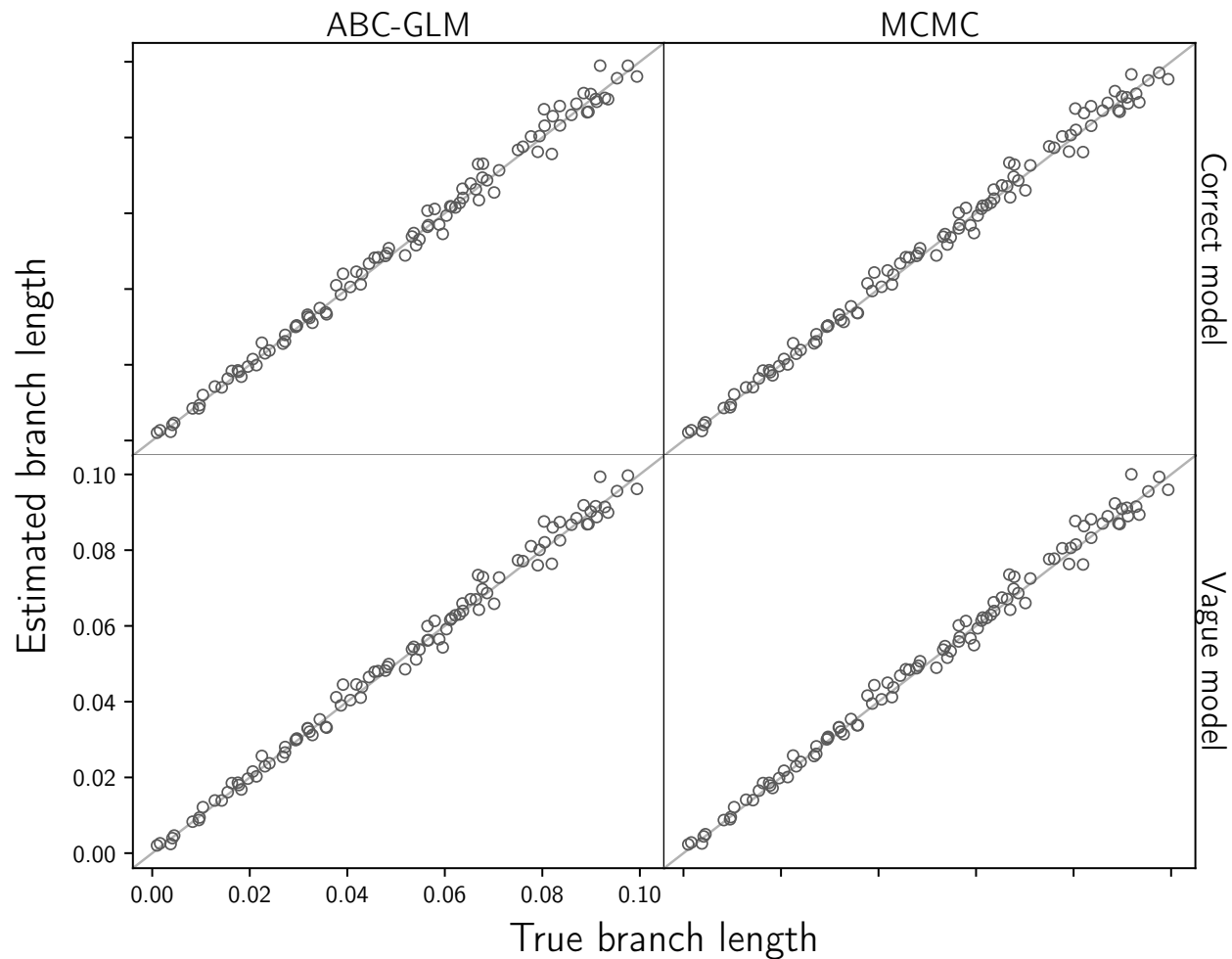


Figure S1. A comparison of the true branch length separating each pair of simulated sequences to the branch length estimated by ABC-GLM and full-likelihood MCMC under the correct branch-length model (branch length \sim uniform(0.0001, 0.1)) and the vague model (branch length \sim uniform(0.0001, 0.1)).