



HAL
open science

Identifying the best approximating model in Bayesian phylogenetics: Bayes factors, cross-validation or wAIC?

Nicolas Lartillot

► **To cite this version:**

Nicolas Lartillot. Identifying the best approximating model in Bayesian phylogenetics: Bayes factors, cross-validation or wAIC?. 2022. hal-04273758v1

HAL Id: hal-04273758

<https://hal.science/hal-04273758v1>

Preprint submitted on 30 Aug 2022 (v1), last revised 7 Nov 2023 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Identifying the best approximating model in Bayesian phylogenetics: Bayes factors, cross-validation or wAIC?

Nicolas Lartillot

Laboratoire de Biométrie et Biologie Evolutive, UMR CNRS 5558, Université Lyon 1, Villeurbanne, France.

`nicolas.lartillot@univ-lyon1.fr`

Running head: Bayes factors, Cross-validation and wAIC

cross-validation, Bayes factor, marginal likelihood, model comparison, wAIC

ABSTRACT

There is still no consensus as to how to select models in Bayesian phylogenetics, and more generally in applied Bayesian statistics. Bayes factors are often presented as the method of choice, yet other approaches have been proposed, such as cross-validation or information criteria. Each of these paradigms raises specific computational challenges, but they also differ in their statistical meaning, being motivated by different objectives: either testing hypotheses or finding the best-approximating model. These alternative goals entail different compromises, and as a result, Bayes factors, cross-validation and information criteria may be valid for addressing different questions. Here, the question of Bayesian model selection is revisited, with a focus on the problem of finding the best-approximating model. Several model selection approaches were re-implemented, numerically assessed and compared: Bayes factors, cross-validation (CV), in its different forms (k-fold or leave-one-out), and the widely applicable information criterion (wAIC), which is asymptotically equivalent to leave-one-out cross validation (LOO-CV). Using a combination of analytical results and empirical and simulation analyses, it is shown that Bayes factors are unduly conservative. In contrast, cross-validation represents a more adequate formalism for selecting the model returning the best approximation of the data-generating process and the most accurate estimates of the parameters of interest. Among alternative CV schemes, LOO-CV and its asymptotic equivalent represented by the wAIC, stand out as the best choices, conceptually and computationally, given that both can be simultaneously computed based on standard MCMC runs under the posterior distribution.

Introduction

Model selection is a difficult question, which has stimulated much theoretical and practical work over the years. The difficulty of the question is due to several factors. First, the models of interest typically differ in their parameterisation, both in structure and in dimensionality, preventing direct comparison of their likelihood scores and requiring careful formalization of how to penalize them accordingly. Second, on a more conceptual front, model selection can be motivated by different objectives, depending on the specific question of interest. These alternative goals entail different compromises and may therefore imply different model selection procedures.

In some cases, the goal of model selection is to test alternative hypotheses about the underlying mechanisms. A relevant example in molecular evolution is the problem of determining whether or not a gene is under positive selection, using phylogenetic codon models. Two alternative models are confronted, one that allows for sites and/or branches under a positive selection regime, tested against a null model that only allows for purifying selection (Nielsen & Yang, 1998; Zhang *et al.*, 2005; Kosakovsky Pond & Frost, 2005). Another example in phylogenetics is the test for the monophyly of a clade. In these examples, the alternative models being considered are meant to be idealized representations of alternative possible states of nature. As a result, the aim is to identify the ‘true’ model, i.e. the model formally representing the true objective situation.

In a classical frequentist context, the standard approach to deal with such hypothesis testing problems is to use likelihood ratio tests, relying on chi-square asymptotics or on parametric (Goldman, 1993) and non-parametric (Shimodaira, 2004) approaches to approximate the distribution under the null. In a Bayesian context, hypothesis testing can be addressed in two different ways. One approach is to compare the marginal likelihoods under the two models, or equivalently, to compute the Bayes factor, i.e. the ratio of the two marginal likelihoods (Jeffreys, 1935; Kass & Raftery, 1995; Oaks *et al.*, 2019). Alternatively, a fully Bayesian formalization of the problem suggests to also define a prior probability over the models and then to select models based on their posterior probabilities (Kass & Raftery, 1995).

In other situations, the question is instead to select the model that gives the most accurate estimation or the best approximation for the data-generating process, and this, without consideration of any hypothesis that would be true or false. A paradigmatic example is to choose the degree of a polynomial regression function (see e.g. Burnham & Anderson, 2002). Here, the true regression function is not generally believed to be itself a polynomial, and thus there is no question

of identifying the true degree. Instead, the question is to find the best tradeoff between the lack of flexibility of polynomials of lower degree and the increased estimation error entailed by a higher degree. Striking the correct balance between these two errors and minimizing the total error is then the fundamental objective of model selection.

In phylogenetics, instances of this second version of model selection are often encountered. An example is the problem of choosing between an empirical matrix such as JTT (Jones *et al.*, 1992), WAG (Whelan & Goldman, 2001) or LG (Le & Gascuel, 2008), or the general time reversible (GTR) model. Empirical matrices are estimates of the average amino-acid exchange rates across a heterogeneous set of proteins and taxonomic groups. As a result, the biochemical prior information that they encode will fit a specific dataset of interest only approximately. If the dataset of interest happens to be sufficiently large, re-estimation of the complete general time-reversible model may give a more accurate model than the one proposed by any available empirical matrix. The problem that model selection has to solve in this context is whether one can afford this re-estimation or whether falling back onto the prior biochemical knowledge encoded into an empirical matrix represents a safer option. The answer to this question will fundamentally depend on data size, but also, on how well the biochemical information encoded into currently available amino-acid replacement matrices generalizes to the specific dataset of interest.

As another example, accounting for pattern heterogeneity across sites is usually done using mixture models (Pagel & Meade, 2004; Koshi & Goldstein, 2001; Lartillot & Philippe, 2004; Quang *et al.*, 2008; Wang *et al.*, 2008; Evans & Sullivan, 2012; Susko *et al.*, 2018; Schrepf *et al.*, 2020). In that context, the question of model selection is important, and non trivial, whether for choosing between alternative empirical models, for determining the number of components, or for the sake of a more general assessment of alternative mixture designs (e.g. finite or infinite mixtures). However, the true distribution of nucleotide or amino-acid substitution rates across sites is not itself a mixture. Instead, the hope is just that a well chosen mixture should give a reasonable approximation of the unknown true distribution, which would then provide increased robustness for phylogenetic inference purposes. The situation is thus formally similar to the one described above in a regression context using a polynomial regression function: the point of model selection with these phylogenetic mixture models is not to identify the true number of components, but to find a good compromise between the lack of flexibility, and potentially the systematic errors, induced by the use of few mixture components, and the increased estimation error entailed by the use of rich mixtures.

The general problem of finding the best approximating model, as opposed to testing hypotheses,

has been classically formalized in different ways. On one side, the approaches used for hypothesis testing, namely likelihood ratio tests and Bayes factors, have often been employed in this context as well. However, it is not totally clear whether they represent a correct formalization of the question, given that there is no proper hypothesis to be tested. As pointed by Akaike (1974) and others (Burnham & Anderson, 2002; Sullivan & Joyce, 2005), hypothesis testing is not adequately formulated, in decision-theoretic terms, as a procedure of approximation, the two goals being intrinsically different. In the more specific context of Bayesian inference, Bayes factors or model posterior probabilities have been recognized as appropriate only in circumstances where it was believed that one of the competing models was in fact true, and that in other circumstances, other criteria may be more appropriate (Bernardo & Smith, 1994; Konishi & Kitagawa, 2007). Accordingly, approaches have been developed, which are more decisively framing the question in terms of finding the best approximation, without predicating on any model being true. Among these approaches, two main types can be identified: cross-validation and information criteria.

The idea of cross-validation is to train the model on a subset of the data and then evaluate the fit of the model over another non-overlapping subset of the observations. The procedure is typically repeated over multiple random splits of the data into a training and a validation set, and the cross-validated log likelihood is finally averaged over these replicates. Cross-validation has been considered both in the context of maximum likelihood (Stone, 1974; Zhang, 1993; Smyth, 2000) and in Bayesian inference (Geisser, 1975; Geisser & Eddy, 1979; Gelfand, 1996). Given its operational definition, cross-validation thus directly estimates the predictive fit of a model. However, this apparent focus on the predictive performance should not be taken too literally. It does not imply that cross validation will be useful only in a context where prediction is indeed contemplated in practice. Perhaps a more fundamental justification is the following: since good prediction of future data can be achieved only by capturing, through the fine-tuning of the parameters of the model, the structural features of the data-generating process, the predictive fit should be good indicator of estimation accuracy. By a similar argument, it can also be seen that cross-validation automatically accounts for overfitting. Indeed, by definition, overfitting is what happens when a model captures random, non-reproducible patterns in the data. Owing to this non-reproducibility, a model that overfits will therefore show a poor fit on new data obtained from the same population. This idea can be quantitatively formalized in terms of the generalization gap of a model (Thomas *et al.*, 2020), or optimism (Efron, 1986), which is defined as the average drop in the apparent log-likelihood score, when going from the training set to the validation set. Altogether, more complex models

will thus have more expressiveness for capturing structural features of the data-generating process, but they will also tend to have a wider generalization gap. Cross-validation automatically captures the balance between these two opposing components of the overall fit.

In the details, cross-validation can be implemented in many different ways, depending on what proportion of the data to set aside for validation, or how many replicates to consider (see Zhang, 1993, for an overview). The simplest and original approach is leave-one-out cross-validation (LOO-CV), whereby each observation is successively taken out of the sample and reserved for subsequent validation of the model, while training is done on the remaining data (Stone, 1974). Alternatively, in k -fold cross validation (k -fold CV), the dataset is split into k equal sized subsets, then each subset is set aside for validation and the remaining $k - 1$ subsets are used for training (Breiman *et al.*, 1984; Zhang, 1993). A variant of k -fold CV is based on repeated random sub-sampling of a fixed fraction $f = 1/k$ of the data that are set aside for validation.

In all cases, direct implementation of cross-validation is expensive, owing to the total number of replicates to consider. Brute-force k -fold cross-validation and its random subsampling version have previously been used in a phylogenetic context (Lartillot *et al.*, 2007; Lartillot & Philippe, 2008), sometimes in combination with strict subsampling, i.e. using training and validation sets that together represent a subset of the data (Pisani *et al.*, 2015). Strict subsampling was motivated by the need to reduce the computational cost. A downside, however, is that the models are then under a regime of data size that does not correspond to the effective regime in which subsequent inference is conducted. Yet the relative fit of alternative models with differing dimensions depends on data size, since higher-dimensional models typically require more data to learn their parameters.

For all these reasons, indirect approaches to cross-validation, which would avoid the explicit resampling and fitting procedure, would be particularly useful. In this direction, and in the specific case of leave-one-out, it is in fact possible to get an estimate of the cross-validation score based only on a standard MCMC run conditioned on the full dataset (Gelfand *et al.*, 1992; Chen *et al.*, 2012; Lewis *et al.*, 2014). This clever importance sampling approach, called cross-predictive ordinates (CPO), makes leave-one-out cross-validation particularly attractive, practically and computationally.

In a more theoretical spirit, and starting with Akaike (1974), a long series of information criteria have been proposed, based on information-theoretic considerations. The fundamental idea behind these information criteria is to identify the model which, once trained on the dataset of interest, induces a distribution over the data that is closest to the true distribution of the population.

Mathematically, the distance between model and truth is measured by the information loss (i.e. the Kullback-Leibler divergence). Importantly, this distance is measured under the effective conditions of use of the model, that is, under the current data size. As a result, it accounts for the two different reasons why the model might not be so close to the true distribution in practice: because of model mis-specification, but also, because of stochastic error in parameter estimation due to finite sample size. This last point will critically depend on both the size of the dataset and the model dimension.

The original criterion proposed by Akaike, the AIC, has a particularly simple expression. However, its derivation relies on the assumption that the models being considered are not far from the true distribution. It is thus not valid under strong model violation, a situation often encountered in practice. The AIC was revisited by Takeuchi (in an original contribution in Japanese, as reported in Konishi & Kitagawa, 1996), who proposed a criterion, the TIC, which is valid even in the presence of strong model violation. The TIC reduces to the AIC when the data are indeed under the model for some true parameter value. Compared to the AIC, the TIC is slightly more involved computationally, since it requires an estimate of the first and second derivatives of the log likelihood at the estimated parameter value. In practice, the difference between TIC and AIC can be substantial (Konishi & Kitagawa, 1996).

The TIC was then generalized, first in a maximum penalized likelihood framework, with the regularized information criterion (RIC, Shibata, 1989) or the generalized information criterion (GIC, Konishi & Kitagawa, 1996) and in Bayesian inference, with the widely applicable information criterion (wAIC Watanabe, 2009). In addition to accommodating model violation like the TIC, the RIC and the GIC, the wAIC is also valid under a broader class of models, such as mixture models or Bayesian networks, which are typically not regular, in the sense that they entail some redundancy (i.e. non-identifiability) in the mapping from parameters to probability distributions over the data (Watanabe, 2007). Because of their non-identifiability, such singular models typically have complex asymptotic properties that are not correctly handled by current information criteria. Addressing these complications is what led to the development of singular statistical learning theory (Watanabe, 2001, 2009), of which the wAIC is one of the specific contributions.

Several other information criteria have been proposed, in addition to those mentioned above. Two of them were explicitly meant for Bayesian inference: the deviance information criterion, or DIC (Spiegelhalter *et al.*, 2002) and the Bayesian analogue of AIC, or AICM (Raftery *et al.*, 2007; Gelman *et al.*, 2014). The DIC has been somewhat controversial (Plummer, 2008; Spiegelhalter *et al.*, 2014; Celeux *et al.*, 2006; Gelman *et al.*, 2014). One problem is that it relies on the posterior

mean point estimate, which is not invariant by re-parameterization of the model and is not easily defined for mixture models (Celeux *et al.*, 2006) or in a phylogenetic context. Another problem is that, like the AIC, the DIC assumes that the model is correctly specified (Spiegelhalter *et al.*, 2002). As for the AICM, it was derived based on an analogy with the AIC, by relying on a definition of the effective number of parameters of a model based on the Monte Carlo variance of the log likelihood. There are two problems with this derivation, however. First, the analogy with the AIC, which is a maximum likelihood criterion, fails to capture the contribution of the prior to the fit of a model in the Bayesian case. Second, just like the AIC and the DIC, the AICM does not account for the impact of model violation. Finally, the Bayesian information criterion, or BIC (Schwarz, 2006) represents one last criterion, which does not proceed from the same rationale as the other criteria mentioned above, as it is not based on an information loss argument. Instead, it is meant as an asymptotic expression for the log of the marginal likelihood. As such it is more appropriate for true model identification than for best model approximation purposes (Aho *et al.*, 2014). Of note, the BIC can be strongly conservative even in a true model identification task (Vrieze, 2012).

There is a direct connection between information criteria based on the information loss (AIC, TIC, RIC, GIC, wAIC) and cross-validation. By definition, the expected information loss is, up to an additive constant that depends only on the true distribution, proportional to the expected log likelihood of new data points sampled from the population, under the parameter value estimated on a data set of the original size. Cross-validation, on the other hand, measures the average log-likelihood of the validation data points, under the parameter value estimated on the training set. Thus, cross-validation can be seen as an operational estimate of the information loss, with the slight nuance that the model is trained on a subset of the data, rather than on the complete dataset. For large datasets, this difference is relatively minor, however, and particularly so for LOO-CV. As a result, LOO-CV is asymptotically equivalent to information criteria of the Akaike family (Stone, 1977; Watanabe, 2010), or equivalently, information criteria of the AIC family are just asymptotic expressions for the leave-one-out cross-validation score, each being valid in a different context and under different specific assumptions. This result is important, as it emphasizes the operational meaning of information criteria. Practically, it suggests simple experiments on empirical data, to check the regime, in terms of data size, under which this asymptotic equivalence is effective (Vehtari *et al.*, 2016) – and thus more fundamentally, the regime under which asymptotic information criteria provide a valid approximation of the frequentist risk that they are intended to measure.

Altogether, there is thus by now a broad theoretical background on model selection. Several

alternative methods have been proposed, with subtle differences concerning their aim or their exact regime of applicability. These issues have already been discussed in the applied statistical literature (Burnham & Anderson, 2002; Aho *et al.*, 2014; Vrieze, 2012; Konishi & Kitagawa, 2007), yet this has not yet been fully incorporated into current phylogenetic practice. This is particularly apparent in Bayesian phylogenetics. Thus, although it has long been noted that Bayes factors are conservative in model selection when used in combination with vague priors on the model specific parameters (the so-called Jeffreys-Lindley paradox, Jeffreys, 1967; Lindley, 1957), and that cross-validation approaches may be more adequate for best-approximating model selection (Gelfand *et al.*, 1992; Bernardo & Smith, 1994; Konishi & Kitagawa, 2007), Bayes factors or marginal likelihoods are often presented as the method of choice (Kass & Raftery, 1995; Lartillot & Philippe, 2006; Xie *et al.*, 2011; Oaks *et al.*, 2019) and are widely used (Suchard *et al.*, 2001; Baele *et al.*, 2012*b*, 2013; Baele & Lemey, 2013; Brown & Thomson, 2017; Ronquist *et al.*, 2021). The computational challenges raised by the numerical evaluation of marginal likelihoods (Lartillot & Philippe, 2006; Xie *et al.*, 2011; Baele *et al.*, 2012*a*) also represent a clear limitation, preventing a broader and more systematic application of this paradigm to current empirical problems based on large datasets. Cross-validation was used in Bayesian phylogenetics primarily for computational reasons (Lartillot *et al.*, 2007; Lartillot & Philippe, 2008), although without any correct evaluation of its numerical accuracy and its theoretical validity in that context. The implementation of LOO-CV offered by CPO appears to be attractive, and has already been introduced specifically in phylogenetics (Lewis *et al.*, 2014), but has thus far not been broadly used in this context. Finally, the wAIC has never been applied to phylogenetic model selection.

In this work, the theoretical and methodological background just presented is utilized to revisit the question of Bayesian model selection in phylogenetics, with an emphasis on identifying the best approximating model, irrespective of any question about hypothesis testing. The statistical and numerical issues are both examined. On the numerical side, the work presented here starts from the realization that k -fold cross-validation, such as implemented in PhyloBayes (Lartillot *et al.*, 2013), turns out to be numerically inaccurate. This point is examined, and an alternative method is proposed, based on sequential importance sampling (sIS), which is similar to sequential Monte Carlo (Wang *et al.*, 2016) and gives an estimate of the marginal likelihood and, simultaneously, the k -fold cross validation scores for any k . This sIS approach is computationally intensive but can be used on datasets of relatively small size to validate and compare marginal likelihood and cross-validation for their ability to select the model that is most accurate in parameter estimation.

Finally, the CPO approach to leave-one-out cross-validation is re-implemented, its statistical and numerical properties are characterized, and its connection with the wAIC is explored on an empirical phylogenomic dataset.

Results

The alternative measures of model fit that are considered in this work are marginal likelihoods, or equivalently Bayes factors, leave-one-out CV (LOO-CV) and k-fold CV (with $k = 5$ and based on independent randomizations of the dataset), the latter in two versions: joint and site-wise. In joint k-fold CV, the joint likelihood of all data points of the validation set is averaged over the posterior distribution under the training set, while in site-wise k-fold CV, the likelihood of each data point of the validation set is averaged over the posterior distribution separately, and then the resulting marginal likelihoods for all data points of the validation set are combined multiplicatively. In order to be measured on the same scale, all scores are log-transformed and normalized so as to be expressed on a per-site basis (see methods for details).

Comparing alternative measures of fit on a simple analytical example

Since they differ in their mathematical definition, these alternative measures of model fit have no reason to agree quantitatively, or even qualitatively, on specific real cases. To examine this point, and before getting into phylogenetic examples, the conceptual and numerical issues are illustrated using simulations under a simple multivariate normal model for which analytical results are available. In this subsection, only the conceptual issues (i.e. the differences in the exact mathematical measures of model fit) are considered, the numerical issues being examined in the next subsection.

The normal model considered here is a variant of the model originally due to Bartlett (1957). The simulated data consist of a series of n real vectors of dimension p , noted $(X_i)_{i=1..n}$, which are iid from a multivariate normal distribution of mean θ_* (also a p -vector) and of covariance matrix $\Sigma = \sigma^2 I_p$, where I_p is the identity matrix. The true mean θ_* used for simulation is chosen to be close to, but not equal to 0. Inference on these simulated data is conducted under two models. In both models, the variance parameter σ is assumed known. Under model M_1 , the vector of means θ is fixed a priori to $\theta_0 = 0$. Under model M_2 , on the other hand, θ is estimated, assuming a normal prior of mean 0 and of covariance $\Sigma_0 = \sigma_0^2 I_p$. The hyper-parameter σ_0 is chosen to be large, so as to implement a vague prior on θ . Of note, when $\sigma_0 \rightarrow \infty$, the prior becomes improper but the

posterior reaches a well-defined limit. We wish to evaluate the relative fit of model M_2 against M_1 on a dataset of size n .

Importantly, the simulation experiment is designed so as to represent a situation where model comparison is recruited for selecting the best approximating model, not the true model. Thus, what we want to formalize is a situation where the fixed parameter value $\theta_0 = 0$ defined by model M_1 is never exactly true. Instead, θ_0 may be viewed as a reasonably good proxy for the unknown true value θ_* , and the question is just whether we can hope to get closer to θ_* by re-estimating θ on the dataset of interest, thus by using M_2 rather than using M_1 .

Data were more specifically simulated under the following settings: $p = 300$, $\sigma^2 = 10$, $\theta_* = 0.1$, and n varying from 100 to 10000. For model M_2 , two values were considered for the prior width, $\sigma_0^2 = 10$ and $\sigma_0^2 = 1000$. Then, the alternative measures of model fit were computed: marginal likelihood (Bayes factor), 5-fold cross-validation, both joint and site-wise, and leave-one-out cross validation. In all cases, the exact analytical values for the expected score of M_2 relative to M_1 were computed. The fit curves are displayed on Figure 1A, as a function of data size. Finally, an analytical formula is also available for the expected root mean squared error under the two models. This expected error, which is thus a frequentist risk, is displayed for the two models on Figure 1B, also as a function of data size.

Several observations can be made from these experiments. First, for small data size, model M_1 is more accurate than model M_2 . Model M_1 is technically wrong (it assumes that $\theta = 0$ whereas in fact $\theta_* > 0$), however, for small data size, the estimation error under model M_2 is much larger than the deviation between θ_* and 0, and thus it is indeed more reasonable to use M_1 in that case. When $n > 1000$, on the other hand, M_2 is more accurate than M_1 .

Second, by comparing the two panels of Figure 1, one can see that Bayes factors are clearly conservative. For instance, when $\sigma_0^2 = 10$, it takes a dataset of at least 8000 observations for Bayes factors to show a preference for M_2 . Thus, between $n = 1000$ and $n = 8000$, Bayes factors are choosing a simple model that can be up to 5 times less accurate than the more complex alternative. This conservativeness is more pronounced under a broader prior (i.e. for larger σ_0). For $\sigma_0^2 = 1000$, the cutoff at which Bayes factors switch to a preference for model M_2 is slightly above $n = 10000$. Importantly, the posterior distribution is virtually the same for these two values of σ_0^2 , which shows that the differences in Bayes factors induced by the choice of the value of σ_0 do not reflect any real-world difference, in terms of estimation. The conservative behavior of Bayes factors under vague priors, such as observed here, is known as Jeffreys-Lindley's paradox (Jeffreys, 1967; Lindley,

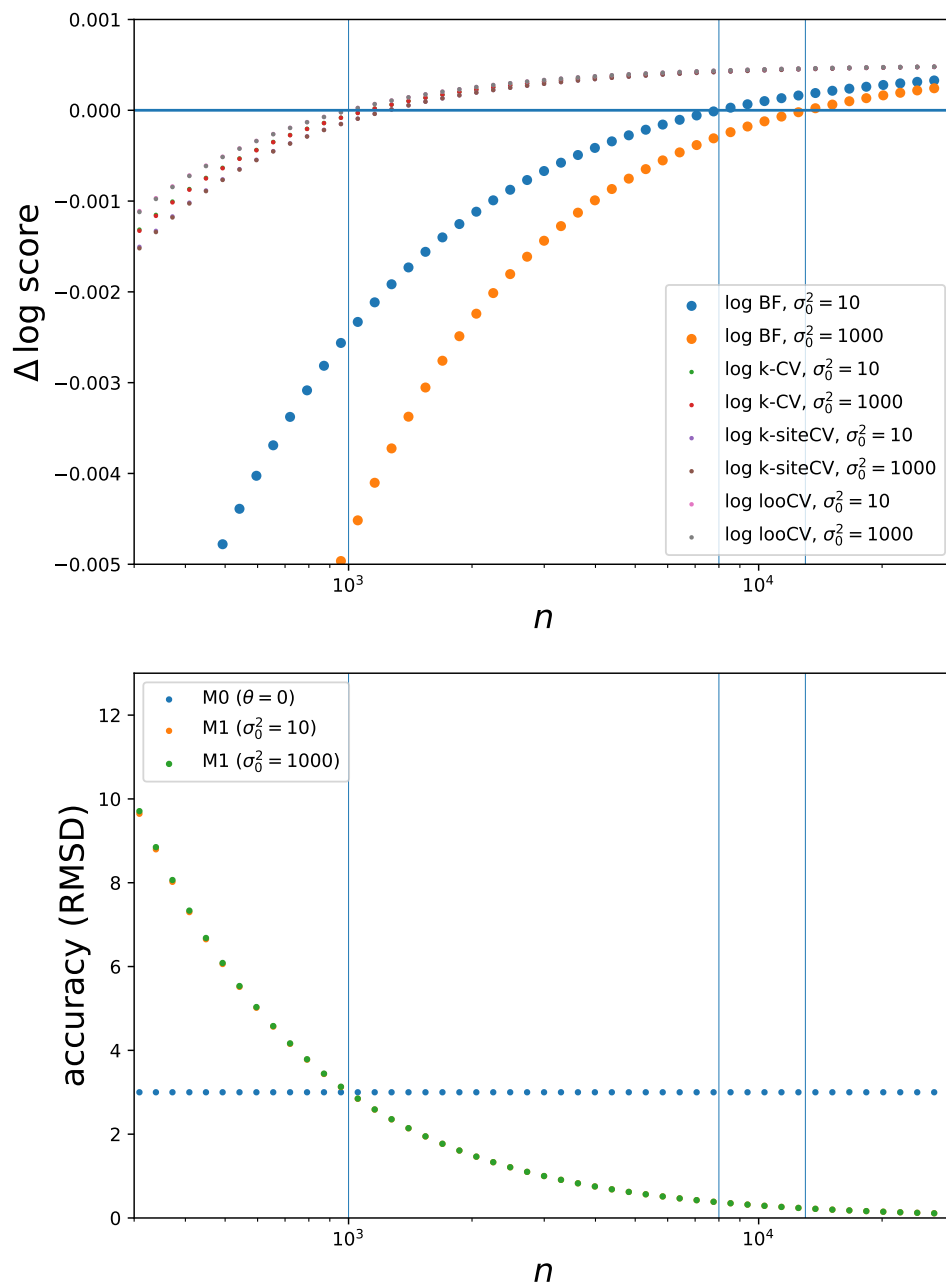


Figure 1: Theoretical fit of M_2 relative to M_1 (top) and mean squared estimation error (bottom) as a function of data size, under the normal model, and for two alternative priors ($\sigma_0^2 = 10$ and 1000); the fit curves under the two alternative priors are indistinguishable for LOO-CV, k-fold joint and site CV.

1957).

In contrast, the model chosen by CV approaches appears to be more directly in proportion to estimation accuracy, with a cutoff very close to the tipping point ($n = 1000$) at which M_2 starts to be more accurate than M_1 . In the details, k-fold CV appears a bit more conservative than LOO-CV, and site-wise k-fold CV is more conservative than both joint k-fold CV and LOO-CV. Although these differences are minor, they illustrate one potential problem with k-fold CV, namely, that it is not measuring the fit under the practically relevant data size. This limitation is inherent to cross-validation, but it is minimized in the case of leave-one-out, for which the training size is virtually indistinguishable from the practically relevant data size for even moderate values of n .

The asymptotic behavior of the alternative measures of fit explored here confirms these points. Up to an order $1/n$, the logarithm of the Bayes factor ($\ln \mathbf{bf}$) and the joint k-fold ($\Delta \mathbf{cv}_j$), site-wise k-fold ($\Delta \mathbf{cv}_s$) and leave-one-out ($\Delta \mathbf{cv}_l$) cv scores of model M_2 relative to M_1 have the following expressions:

$$\ln \mathbf{bf} \simeq \frac{p \theta_*^2}{2 \sigma^2} - \frac{p \ln n}{2 n} + \frac{p}{2n} \left(1 - \frac{\theta_*^2}{\sigma^2} - 2 \ln \left(\frac{\sigma_0^2}{\sigma^2} \right) \right) \quad (1)$$

$$\Delta \mathbf{cv}_s \simeq \frac{p \theta_*^2}{2 \sigma^2} - \frac{p}{2(1-f)n} \quad (2)$$

$$\Delta \mathbf{cv}_j \simeq \frac{p \theta_*^2}{2 \sigma^2} - \frac{p}{2n} \frac{|\ln(1-f)|}{f} \quad (3)$$

$$\Delta \mathbf{cv}_l \simeq \frac{p \theta_*^2}{2 \sigma^2} - \frac{p}{2n} \quad (4)$$

Of note, when the set-aside fraction f is small, then $\frac{1}{1-f} \simeq 1 + f$, and $|\ln(1-f)| \simeq f + \frac{1}{2}f^2$, such that:

$$\Delta \mathbf{cv}_s \simeq \frac{p \theta_*^2}{2 \sigma^2} - \frac{p}{2n} (1 + f) \quad (5)$$

$$\Delta \mathbf{cv}_j \simeq \frac{p \theta_*^2}{2 \sigma^2} - \frac{p}{2n} \left(1 + \frac{f}{2} \right) \quad (6)$$

As for the asymptotic relative risk (i.e. difference in quadratic error between model M_1 and M_2 , normalized here by $2\sigma^2$), it is, up to terms in $1/n$:

$$\Delta R = \frac{p \theta_*^2}{2 \sigma^2} - \frac{p}{2n} \quad (7)$$

From these equations, several observations can be made. On one hand, for sufficiently large n , all terms except the first vanish. Since this first term is positive, all measures will eventually agree and will all choose M_2 . This asymptotic agreement is visible on Figure 1A. Also visible on the figure is

the slower convergence, in $\ln n/n$, for the log Bayes factor, whereas it is in $1/n$ for cross-validation under all settings. Of note, this mirrors the penalties of the BIC and the AIC, respectively.

Conversely, however, for fixed n , and considering increasingly vague priors by letting σ_0 go to infinity, the log Bayes factor is ill-behaved, since its last term goes to $-\infty$. In other words, for a given data size, and for arbitrary large θ_* , then, provided that the prior is sufficiently broad, BF will nevertheless prefer M_1 – and this, in spite of the arbitrary large risk that this might entail. In contrast, CV measures are all well-behaved and are insensitive to σ_0 . When the set-aside fraction f is small, the two versions of k -fold CV are slightly more conservative than LOO-CV (that is, they have a slightly stronger penalty), the joint version being intermediate between LOO-CV and the site-wise version, as seen on Figure 1A. Finally, LOO-CV is asymptotically equal to the difference in quadratic estimation error between the two models. In other words, asymptotically, LOO-CV is exactly selecting the model that gives the most accurate estimate. However, this last point not a general result. Instead, it is a consequence of the fact that a spherical covariance matrix was used in the model. For general covariance structures, the loo-score is asymptotically equal to another relative quadratic risk, computed under the metric defined by the covariance matrix. This metric essentially gives less weight to the errors made on those components of θ for which the likelihood is less informative.

A final point not quantitatively explored here but worth noting: when $\theta_* = 0$, that is, when M_1 is the true model, all CV methods considered here are asymptotically inconsistent, in the sense that the probability of choosing M_1 does not converge to 1 for large n (Shao, 1993). However, whenever CV chooses M_2 , it will then estimate a value for θ very close to its true value 0 (up to a quadratic error in $1/n$), such that the error in model selection will have a negligible impact on estimation accuracy. In other words, CV is not formally consistent, but it is effectively consistent, in the sense that the selected model is asymptotically equivalent to the true model in the Kullback-Leibler metric. Conversely, since LOO-CV is asymptotically optimal in estimation accuracy in the present case, any method trying to be asymptotically formally consistent will have to be more strongly penalizing than LOO-CV and will thus be suboptimal for selecting the best approximating model. The two goals of model selection, best approximation or true model identification, are thus mutually incompatible (Shibata, 1986).

Numerical approaches: accuracy and computational complexity

In this section, the question of the numerical evaluation of the alternative measures of model fit considered above is explored, again in the normal case, for which the numerical estimates can be directly assessed against the analytically available value. The Monte Carlo approaches that are used here are all variations on importance sampling: naive importance sampling (nIS) for both joint and site-wise k-fold CV, sequential importance sampling (sIS) for joint k-fold CV and marginal likelihoods, and the cross-predictive ordinate (CPO) approach for LOO-CV (see methods for details).

Naive importance sampling simply consists of averaging the likelihood of the data points of the validation set (either jointly or separately, for joint or site-wise k-fold CV respectively) over a sample of parameter configurations drawn from the posterior distribution under the training set. When applied to joint k-fold CV, nIS works well for low dimension ($p = 10, 30$ or 100) but its performance progressively degrades as the dimension of the model increases. For large dimension $p > 300$, a substantial downward bias is observed. In the case of $p = 1000$, the bias is sufficiently strong to change the qualitative outcome of model selection, leading to an apparent CV score in favor of model M_1 , whereas model M_2 has mathematically a higher CV score. As expected, increasing the size of the Monte Carlo sample can improve the situation, although very moderately. Under the highest dimensions considered here, it seems that it would take samples of very large size, well above 10^6 , in order to reduce the bias down to reasonable values.

A key statistic that is able to issue a warning about the reliability of the estimation in the present case is the effective sample size (ESS). The ESS is a function of the variance of the importance weights, such that the ESS is close to 1 when a single point of the sample has an overwhelming contribution to the Monte Carlo average (essentially, the point of the sample that happens to have the highest likelihood). In the present case, for high dimensions, the ESS is indeed close to 1, indicating that the estimator is fundamentally unreliable.

In contrast to what is observed for joint k-fold CV, nIS works well on site-wise k-fold CV (Table 1). This is due to the fact that the single-observation likelihood is much less peaked than the joint likelihood of multiple observations. As a result, the variance of the log-likelihood score under the posterior distribution is small. Of note, for small MCMC sample size (10 samples per site), the total bias of the estimator in log scale can be non-negligible. On the other hand, because this bias is a sum of many small contributions (one for each data point), each of which has a large ESS

method	dim	sample size		model fit			bias		error	
		nominal	ESS	true	est. ¹	deb. ²	true	est.	raw ³	deb. ⁴
k-CV (nIS)	100	10 ⁴	6.49	0.02	0.017	0.018	-0.006	-0.001	0.009	0.009
	100	10 ⁶	24.20	0.02	0.021	0.021	-0.002	-0.000	0.004	0.004
	300	10 ⁴	2.26	0.08	0.010	0.011	-0.065	-0.001	0.067	0.066
	300	10 ⁶	2.95	0.08	0.033	0.035	-0.042	-0.001	0.043	0.042
	1000	10 ⁴	1.52	0.24	-0.159	-0.157	-0.394	-0.002	0.395	0.393
	1000	10 ⁶	1.74	0.24	-0.091	-0.090	-0.327	-0.002	0.328	0.326
k-site-CV (nIS)	100	10	9.01	0.02	0.008	0.015	-0.007	-0.006	0.012	0.010
	100	10 ³	883.02	0.02	0.015	0.015	-0.000	-0.000	0.001	0.001
	300	10	7.55	0.06	0.035	0.055	-0.020	-0.020	0.026	0.017
	300	10 ³	689.11	0.06	0.055	0.056	-0.000	-0.000	0.002	0.002
	1000	10	5.15	0.16	0.063	0.127	-0.101	-0.064	0.106	0.050
	1000	10 ³	302.75	0.16	0.162	0.163	-0.002	-0.001	0.004	0.004
k-CV (sIS)	100	10	9.09	0.02	0.018	0.023	-0.006	-0.006	0.008	0.006
	100	10 ³	894.89	0.02	0.023	0.023	-0.000	-0.000	0.001	0.001
	300	10	7.74	0.08	0.057	0.077	-0.018	-0.020	0.020	0.010
	300	10 ³	716.99	0.08	0.075	0.075	-0.000	-0.000	0.001	0.001
	1000	10	5.39	0.24	0.153	0.236	-0.082	-0.083	0.084	0.023
	1000	10 ³	341.29	0.24	0.234	0.235	-0.001	-0.001	0.003	0.002
BF (sIS)	100	10	8.91	-0.33	-0.339	-0.331	-0.008	-0.007	0.008	0.003
	100	10 ³	871.21	-0.33	-0.331	-0.331	-0.000	-0.000	0.000	0.000
	300	10	7.39	-1.00	-1.020	-0.995	-0.024	-0.025	0.024	0.005
	300	10 ³	663.74	-1.00	-0.996	-0.996	-0.000	-0.000	0.001	0.001
	1000	10	4.98	-3.30	-3.420	-3.310	-0.116	-0.110	0.117	0.013
	1000	10 ³	277.90	-3.30	-3.306	-3.304	-0.002	-0.002	0.002	0.001
LOO-CV (CPO)	100	10	9.18	0.03	0.035	0.030	0.005	0.005	0.006	0.003
	100	10 ³	904.89	0.03	0.030	0.030	0.000	0.000	0.000	0.000
	300	10	7.93	0.09	0.103	0.086	0.018	0.017	0.019	0.006
	300	10 ³	741.22	0.09	0.085	0.085	0.000	0.000	0.001	0.001
	1000	10	5.60	0.30	0.374	0.302	0.073	0.072	0.074	0.012
	1000	10 ³	376.91	0.30	0.302	0.301	0.001	0.001	0.001	0.001

¹: estimated; ²: debiased; ³: true error of raw estimator; ⁴: true error of debiased estimator

Table 1: Numerical estimates of the fit of M_2 (relative to M_1) under various criteria and numerical approaches for the normal model example. See text for details.

and therefore a small variance, it can be estimated based on a linear approximation relating it to the variance observed across two independent runs (see methods). As a result, it can be worth de-biasing the estimator; although not perfect, doing this does increase the overall accuracy, quite substantially (Table 1).

A sequential importance sampling approach for k-fold CV and BF

As a way to overcome the limitations of naive IS, an alternative approach was implemented, based on sequential importance sampling (sIS). Briefly, the idea of sIS is to run a quasi-static MCMC in which data points are introduced one by one, each time running the MCMC for a few cycles for equilibration, followed by another series of cycles for averaging the likelihood of the next data point under the posterior induced by all current data points. This can be seen as a variation on stepping-stone or on thermodynamic integration (Lartillot & Philippe, 2006; Fan *et al.*, 2011; Xie *et al.*, 2011), in which the power posteriors have been replaced by partial posteriors (i.e. based on increasingly large subsets of the data). As such, it is also close in spirit (although slightly different from) sequential Monte Carlo (Wang *et al.*, 2016).

When applied to the multivariate normal problem, sIS gives a more reliable estimate of the joint k-fold CV score over the whole range of model dimensionalities considered here (Table 1). The estimate of the log marginal likelihood returned by sIS is also reliable, for both small and large dimensions. Here again, as in the site-wise case, the total bias of the estimators can be substantial for small Monte Carlo sample size, but is itself well-estimated. This sIS approach, however, is expensive – even more expensive for CV than for marginal likelihood, since CV requires to run sIS ideally over a large number of randomized replicates of the original dataset, whereas only two runs on the original non-randomized alignment are needed for the marginal likelihood.

Leave-one-out cross validation using cross-predictive ordinates

An estimate of the LOO-CV score can be obtained very efficiently, based on a standard MCMC run under the posterior distribution, using the CPO approach (Gelfand *et al.*, 1992; Chen *et al.*, 2012; Lewis *et al.*, 2014). The CPO method gives accurate estimates of the LOO-CV score (Table 1). Here again, the bias can be substantial for small sample size but is well estimated.

Altogether, naive IS works well for site-wise k-fold CV, but does not work well for joint k-fold CV. Both joint k-fold CV and Bayes factors require computationally more intensive MCMC approaches, such as the sequential IS approach used here. Finally, the CPO approach represents a

model	5-fold CV						BF		LOO-CV	
	nIS (10^3)		nIS (10^4)		sIS		sIS		CPO	
	fit	(ESS)	fit	(ESS)	fit	(ESS)	fit	(ESS)	fit	(ESS)
Poisson	-30.89	(3.3)	-30.88	(6.7)	-30.87	(28.4)	-31.49	(28.3)	-31.35	(856.5)
WAG	-28.18	(5.7)	-28.17	(8.3)	-28.16	(28.3)	-28.82	(28.2)	-28.68	(864.2)
LG	-27.96	(5.1)	-27.95	(6.2)	-27.94	(28.3)	-28.61	(28.3)	-28.47	(864.5)
GTR	-27.92	(1.8)	-27.89	(2.3)	-27.79	(31.6)	-28.65	(31.8)	-28.31	(675.2)
CAT-fix-hyper	-27.54	(1.8)	-27.51	(2.5)	-27.34	(33.6)	-28.10	(36.8)	-27.80	(658.1)
CAT-free-hyper	-27.66	(1.3)	-27.59	(1.5)	-27.18	(37.1)	-28.04	(36.0)	-27.67	(552.0)

Table 2: Numerical estimates of the fit of alternative substitution models for the elongation factor alignment. All measures of fit are on a logarithmic scale and on a per-site basis.

reliable and computationally efficient method for estimating the LOO-CV score.

An empirical example using a single-gene alignment

The various scores and numerical methods for computing them, such as explored above on the normal case, were then implemented in PhyloBayes (Lartillot *et al.*, 2009, 2013). In a phylogenetic context, it is natural to use the individual columns of the multiple sequence alignment as the individual data points. For the rest, the implementation of all of the methods is relatively straightforward, based on the already existing MCMC routines. All of these estimators were then jointly examined, in the context of a global comparison between alternative site-homogeneous and site-heterogeneous models on an empirical alignment. The models under comparisons are the Poisson model (exchangeabilities between amino-acids all equal to 1), the empirical matrices WAG (Whelan & Goldman, 2001) and LG (Le & Gascuel, 2008), and finally, the CAT-Poisson model (Lartillot & Philippe, 2004), in two alternative versions that differ in the base distribution used for the Dirichlet process over the amino-acid frequency vectors: either a uniform (fix-hyper) or a general (free-hyper) Dirichlet distribution whose hyperparameters are then also estimated. The latter is the version of the CAT model proposed by default by PhyloBayes. The results are presented in Table 2.

First, concerning k-fold CV, nIS and sIS (which are two alternative estimators of the same mathematical quantity) agree with each other on simple models such as Poisson or WAG, but not for more complex models such as CAT-Poisson. The ESS clearly suggests that, here also, as in the normal case, nIS is being unreliable. In one case, this leads to a different qualitative answer as to which model is best fitting. Thus, nIS gives an apparently higher joint k-fold CV score for the

version of the CAT model that uses a uniform base distribution (fix-hyper), whereas sIS says that the version with a general Dirichlet distribution (free-hyper) has a higher fit. Of note, fix-hyper is a constrained version of the free-hyper model, assuming a uniform base distribution (as opposed to a hyper-parameterized Dirichlet distribution). The posterior estimate of the base distribution under the unconstrained model, however, is very far from uniform, which gives an additional argument, independent from the ESS, suggesting that nIS is giving a wrong answer in the present case.

Second, concerning the alternative measures of model fit: k-fold CV and LOO-CV qualitatively agree with each other. On the other hand, they differ somewhat from the Bayes factor, which tends to be more conservative. In one case, BF and CV give a qualitatively different outcome, concerning the choice between GTR and empirical matrices: whereas all CV methods choose GTR, BF gives a higher score to the LG model on this EF2 dataset.

LOO-CV, BF and estimation accuracy in a phylogenetic context

The experiment above on EF2 suggests that BF can sometimes disagree with cross-validation on real cases. To further investigate this point, another experiment was conducted, consisting of comparing LG and GTR on increasingly large subsets of an empirical supermatrix of 35 metazoan species (Philippe *et al.*, 2005), using either BF or LOO-CV. The k-fold CV approach was not considered, owing to its computational cost. Of note, here as above (Table 2), the prior on the renormalized exchangeabilities (constrained to sum to 1) of the GTR model is uniform, thus uninformative.

The results of this experiment are summarized in Figure 2. For sufficiently large datasets, BF and LOO-CV both favor the GTR model over LG, while for smaller datasets, the LG model tends to be favored. This point is expected, and confirms that, for sufficiently large data size, there is an opportunity for getting better estimates of the relative exchangeabilities than those proposed by LG. However, if both methods agree on this dichotomy between small versus large datasets, they substantially differ concerning the exact cutoff, in terms of data size, at which they switch from LG to GTR. Here again, as seen above in the case of the normal model, BF is generally more conservative than LOO-CV. As a consequence, there is an intermediate regime for which BF and CV qualitatively differ in their selection. In practice, this intermediate regime covers a non-negligible interval: whereas BF favors GTR over LG only starting from alignments made of more than 600 sites, LOO-CV does so for datasets as small as 200 sites.

The analytical results presented above under the normal model suggested that cross-validation is more in phase with estimation accuracy than Bayes factors. To investigate whether this conclusion

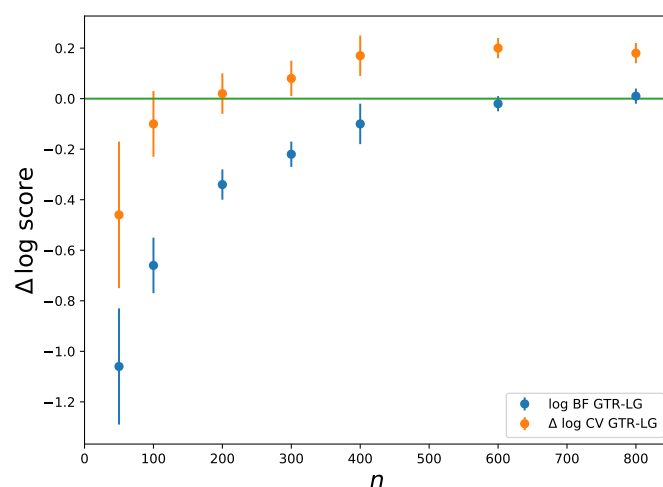


Figure 2: Fit of GTR, relative to LG, as a function of data size (number of aligned positions), on empirical data (10 random jackknife subsamples of the metazoan dataset), using BF and LOO-CV. Error bars: standard deviation across jackknife replicates.

is also valid in the present case, the following simulation experiment was conducted. First, data were simulated under LG and under empirically calibrated branch lengths and parameter values, using the posterior predictive formalism and with the metazoan alignment as a template (see methods). Then, model selection was implemented, between JTT and GTR. Importantly, the true model (LG) was not included in the set of models being compared. This omission is meant to represent the fact that, in reality, the true exchange rates (or, more accurately, the asymptotic exchange rates, i.e. the ones that would be eventually estimated on a sufficiently large alignment obtained from this empirical source) are not equal to any of the empirical models that are available. The difference between JTT and LG is thus meant as a representation, in our simulation experiment, of the difference between LG and the true exchange rates in the empirical experiment.

The Bayes factors and cross-validation scores obtained on these simulated data (Figure 3a) reproduce the pattern observed on the empirical data (Figure 2) as a function of data size, with GTR being ultimately favored by both BF and LOO-CV, although for a larger cutoff data size for BF (700) than for LOO-CV (200). Of note, both the cutoffs and the absolute fit values are very similar to those obtained on the original empirical experiment (Figure 2), suggesting that the simulation experiment is mimicking the true empirical situation relatively well.

Along with model fit, the error (RMSD) in the estimation of the relative exchange rates was

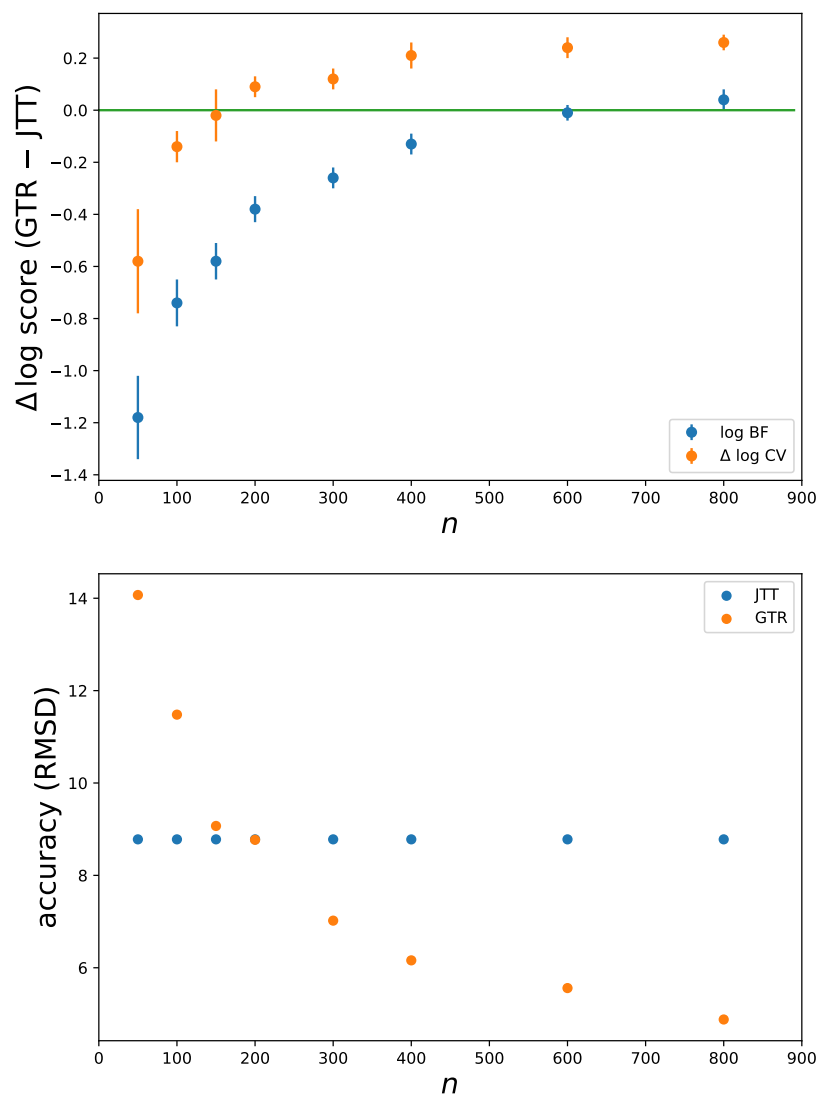


Figure 3: Fit of GTR, relative to JTT, as a function of data size under BF and LOO-CV (top), and mean quadratic error on relative exchangeability estimates (bottom) on data simulated under the LG model (using the metazoan dataset as a template). Error bars: standard deviation across 4 simulation replicates.

also quantified (Figure 3b). In the case of the JTT model, this error is trivially constant (quadratic deviation between JTT and LG). For the GTR model, this error decreases with data size. On sufficiently small alignments, on the other hand (smaller than 200), the estimation error under GTR can be larger than the difference between JTT and the true exchange rates (LG). Thus, for small alignments, we are in a case where the most accurate model is in fact JTT, and this, in spite of the fact that JTT is not the true model.

Finally, comparing RMSD with both BF and LOO-CV shows that LOO-CV provides a good predictor of which model is more accurate for parameter estimation, with a cutoff at around 200 aligned positions. BF in contrast, imposes a stronger penalty and still chooses JTT for datasets up to 600 sites, thus well within the regime of alignment size where GTR is in fact already returning a substantially more accurate estimation. Transposing these observations to the empirical case, this suggests that LG is in fact not so good and GTR is better, even for small alignments of about 200 sites and 50 taxa. It also confirms the point already demonstrated on the normal case, namely that LOO-CV gives a more reliable predictor of estimation accuracy than BF.

Asymptotics of LOO-CV and the widely applicable information criterion (wAIC)

Bayesian LOO-CV is asymptotically equivalent to the wAIC, which is an adaptation of the AIC and, more fundamentally, the TIC, to the Bayesian case (Watanabe, 2007). The wAIC (per site) takes the following form :

$$\text{wAIC} = \frac{1}{n} \sum_i \ln E_{post}[p(X_i | \theta)] - \frac{1}{n} \sum_i V_{post}[\ln p(X_i | \theta)] \quad (8)$$

where E_{post} is the expectation, and V_{post} the variance, over the posterior distribution under the complete dataset X . In practice, these theoretical expectation and variance terms are replaced by their Monte Carlo counterparts (empirical mean and variance over the MCMC sample).

In terms of interpretation, the first term of wAIC can be seen as the self-fit, that is, the fit of the training set under the parameter value estimated on that training set. Because it uses the data twice, this measure of the fit is optimistic. The second term represents an estimate of this optimism bias. As such, it plays the same role as the dimensional penalty in AIC. Of note, in spite of their similar form, the two terms in equation 8 are not of the same order of magnitude. Owing to the asymptotic concentration of the posterior, the variance of the log likelihood at a typical site (the second term) decreases as a function of data size, whereas the mean log likelihood (the first term) remains asymptotically macroscopic. The same situation holds for the AIC and other classical

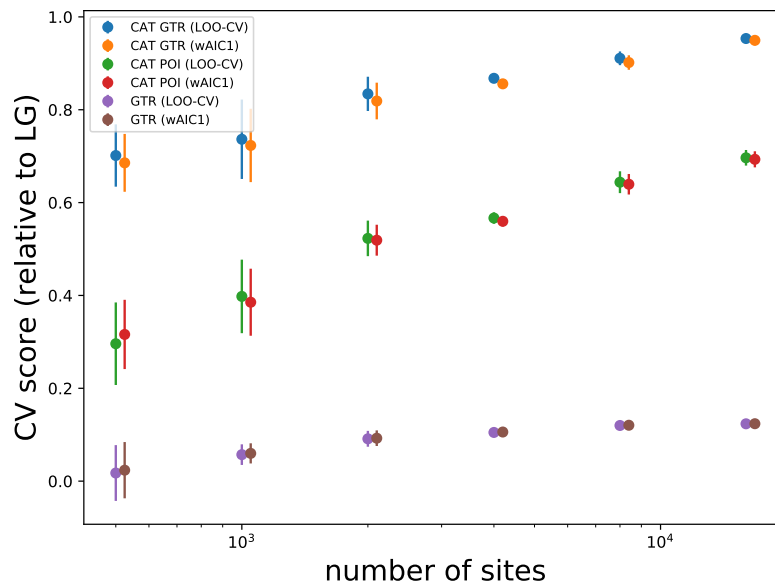


Figure 4: LOO-CV and wAIC for the GTR, CAT-Poisson and CAT-GTR models (relative to LG), as a function of data size (number of aligned positions), on empirical data (metazoan dataset). Error bars: standard deviation across 4 jackknife replicates.

information criteria, for which the dimensional penalty becomes negligibly small compared to the log likelihood term for sufficiently large datasets.

The numerical stability and the asymptotic equivalence between the wAIC and LOO-CV (Watanabe, 2010) were empirically assessed by conducting a scaling experiment, consisting of randomly subsampling a large phylogenomic dataset and plotting the fit (LOO-CV and wAIC) of the GTR, CAT-Poisson and CAT-GTR models (relative to LG), as a function of data size (Figure 4). Overall, LOO-CV and the wAIC give very similar results. The discrepancy between them decreases as the data size becomes larger, giving nearly indistinguishable numerical estimates for the largest data sizes considered here. Even for smaller data size, the difference between LOO-CV and the wAIC is visible but small compared to the difference in fit between the models.

In terms of numerical stability, both estimators, of LOO-CV and of wAIC, tend to be numerically more stable for larger data size. In addition, the wAIC tends to be more accurate and less sensitive to variation in Monte Carlo sample size than LOO-CV (Supplementary Information, section 2). In the end, for large datasets (more than 4000 positions), for which the asymptotic approximation of the wAIC is very accurate, it is possible to get numerically satisfactory estimates of the wAIC based

on a thinned sample of 100 points regularly spaced over a well-converged MCMC run. This thinning approach represents a particularly attractive option in the case of mixture models, for which the numerical evaluation of the likelihood, as a sum over mixture components, is computationally intensive and turns out to be the limiting factor for estimating the fit.

A final note concerning the learning curves of the three models: for the GTR model, the learning curve reaches a plateau at around 5000 aligned positions. This suggests that, with 5 000 sites, the GTR model has essentially learned all of what it could learn from the data. For the CAT-Poisson and CAT-GTR models, on the other hand, the learning curves do not appear to stabilize, even for the largest alignments (16 000 positions). Thus, although CAT-Poisson and CAT-GTR have a substantially higher fit than GTR over the whole range of data size considered here, they could apparently do much better still with more data or, in other words, there is still much to learn about the details of the distribution of amino-acid preferences across sites.

Discussion

In many respects, model comparison and model selection in Bayesian inference is still an open problem. Conceptually, in spite of a large literature on the question, a general agreement on the guiding principles has not yet been achieved. Computationally, numerical inaccuracies are surfacing regularly. The present work attempts to bring a few points of clarification, along with a correction concerning the numerical accuracy of a previously introduced importance sampling k-fold CV approach. In the end, some recommendations are suggested for improving both reliability in model selection and computational accuracy and efficiency.

The main conclusions are as follows. As suggested previously (Gelfand *et al.*, 1992; Bernardo & Smith, 1994; Konishi & Kitagawa, 2007), Bayes factors are inadequate for selecting the best-approximating model, and cross-validation appears to be more adequate for this purpose. Among CV methods, LOO-CV stands out as the best choice, both statistically and computationally. It also has a clear asymptotic connection with information criteria, and more specifically with the widely applicable (or Watanabe-Akaike) information criterion (wAIC Watanabe, 2009). For large datasets, wAIC is easily implemented and offers a good complement to LOO-CV.

Problems with marginal likelihoods under vague priors

One first fundamental reason for the conservative behavior of marginal likelihoods in the present case is the use of a vague prior over the model-specific parameters. From the experiments presented here, and more generally on conceptual grounds, marginal likelihoods do not represent a meaningful measure of model fit under a prior that is meant to be uninformative. This is particularly apparent in the case of the normal model. Under this model, when the prior over the unknown mean θ is uniform over the entire real line, and thus improper, the posterior distribution is well-defined, but the marginal likelihood is infinite. This problem has been known for a long time (Gelfand *et al.*, 1992; Jeffreys, 1967; Lindley, 1957), and it has already been noted in the context of phylogenetic inference that marginal likelihoods and Bayes factors should not be used with improper priors (Baele *et al.*, 2012b). However, making the prior technically proper but still effectively uninformative does not solve the problem. This is again clear in the case of the normal model, for which model selection based on the marginal likelihood can be made arbitrarily stringent against the more complex model by playing on its width parameter σ_0 – and this, in spite of the fact that the posterior distribution is virtually unaffected (Figure 1b, compare red and green dots). This problem is also well illustrated by the comparison between JTT and GTR. In that case, the prior over the relative exchangeabilities of the GTR model is proper, but non informative, and the marginal likelihood is unduly biased in favor of JTT.

A reasonable operational consistency requirement in the context of best-approximating model selection would be that the criterion used for selecting models should give essentially identical scores to models that give essentially identical posterior distributions. Obviously, marginal likelihoods do not fulfill this consistency requirement. They are notoriously sensitive to the prior – and more so than the posterior distribution itself. In contrast, cross-validation, and its asymptotic equivalent given by information criteria such as the wAIC, are by construction dependent on the prior only through the posterior distribution. Thus, they are guaranteed to be operationally consistent.

Importantly, all this does not imply that using uninformative priors is in itself problematic. Uninformative priors do have a good theoretical justification, as a bet-hedging strategy, whose aim is to minimize the worst case error over all possible values the unknown parameter might have (Berger, 1985). As such, they are generally proposed as default priors, meant to guarantee some robustness in the context of automatic application of the inference method to an arbitrary series of practical cases (Berger, 2006). They are thus particularly useful as routine priors, in particular for

the global parameters of the model, which are at the top of the hierarchy. However, model selection methods should then be compatible with these priors.

Cost of learning versus accuracy, and the two aims of model selection

Intuitively, another fundamental problem of marginal likelihoods in the present context is that they penalize models in proportion to how much information has been extracted from the data (essentially, the relative width of the posterior, compared to the prior), and not in proportion to how accurately this information has been learned. Yet, for selecting the best-approximating model, only the second point is relevant. The penalty induced by cross-validation, on the other hand, is directly and exclusively related to how well the fitted model predicts new data. As a result, it is more directly related to the accuracy of the end result of the estimation, leaving out any consideration about the total cost of parameter fitting.

The sequential importance sampling formalism gives another intuition of the same idea. With sIS, the logarithm of the marginal likelihood is obtained by starting from the prior, adding sites incrementally and summing up their individual contributions. Thus, the overall score is a sum over the total learning curve, and as a result, it penalizes models in proportion to the total learning work done upon going from the prior to the posterior. In contrast, leave-one-out cross validation considers only the last step of the procedure and therefore penalizes in proportion to the marginal surprise of the last data point (taken as a proxy for the average future observation). Thus, again, cross-validation is more directly related to the operational quality of the final outcome, not to the entire process of model fitting.

This difference between the two approaches is reflected in the scaling of the asymptotic penalties that were derived in the case of the normal model (equations 1 to 4) but that are more generally valid: on a per-site basis, the penalty is in $\ln n/n$ for the marginal likelihood (or the BIC), and in $1/n$ for cross-validation (or the AIC and its relatives). In turn, $\ln n$ is an asymptotic for the cumulative sum of $1/k$:

$$\sum_{k=1}^n \frac{1}{k} \sim \ln n$$

which thus reveals that the penalty of BF and BIC is indeed capturing the total cost of fitting (as if sites had been added one by one), as opposed to the marginal cost of fitting of the last observation for CV and AIC.

On the other hand, the total cost of fitting and, more generally, the sensitivity of the marginal

likelihood to the prior, is potentially relevant for hypothesis testing. For instance, an alternative hypothesis may explain the data better than does the null, but only under an effect size that is very small compared to the typical effect sizes that would be a priori expected if the alternative were true. This a priori unlikely event will represent a cost that marginal likelihoods will incorporate in their evaluation of the fit, making them more inclined to select the null hypothesis in that case. Marginal likelihoods will thus be useful in a hypothesis testing context, although this requires careful design of the priors over effect sizes and over alternative hypotheses, so as to ensure a correct calibration of the test. This point, and more generally the question of Bayesian hypothesis testing and its application to phylogenetic problems, certainly deserve further investigation.

In contrast, cross-validation, being insensitive to the cost mediated by the prior on the effect size, will often incorrectly choose the alternative in this hypothesis testing example. More generally, cross-validation will often fail at suppressing minor but irrelevant fluctuations and redundancies from the output and, as a result, will not be asymptotically consistent in true model identification (Shao, 1993). However, this may be the price to pay, in order to obtain a model selection criterion that is sufficiently flexible in other contexts and for other purposes, such as fitting a sufficiently fine-grained mixture to a complex distribution of random effects. The different aims of model selection, testing hypotheses or finding the best-approximating model, just entail different compromises (Shibata, 1986).

Implications for Bayesian model averaging

Model averaging is a powerful feature of Bayesian inference, making it possible to consider large combinatorial spaces of model configurations, while integrating uncertainty over models, effect sizes and nuisances (Fragoso *et al.*, 2017; Hoeting *et al.*, 2000). However, Bayesian model averaging implicitly relies on marginal likelihoods. Therefore, when used in combination with uninformative priors, it will also be biased in favor of the simpler models, just like marginal likelihoods and Bayes factors in the context of explicit model selection. This potentially concerns several previously introduced approaches, implementing model averaging over nucleotide substitution models (Huelsenbeck *et al.*, 2004), over the number of components of a mixture (Evans & Sullivan, 2012), or over the number of change points of a non-homogeneous substitution model along the phylogeny (Blanquart & Lartillot, 2006). In the cases just cited, an uninformative prior is used, not just for the global parameters of the model, but also for the replicated items (the exchange rates, the mixture components or the effect sizes associated with each change point). As a result, there is a

tendency to over-penalize the more complex model configurations, in spite of the fact that those might be empirically more adequate.

Over-penalization in the context of Bayesian model averaging can be mitigated by the use of a hierarchical prior over the replicated items. For instance, in the case of mixture models, using a hyper-parameterized prior over component-specific parameters will make the model averaging approach less conservative and thus empirically better fitting – as can be seen when comparing the hierarchical (free-hyper) version of the CAT model with its non-hierarchical (fix-hyper) version (Table 2). Similarly, under the change-point model, hyperparameterizing the distribution of the effect sizes upon each transition will result in a more flexible and empirically more adequate model.

All of these points certainly need further exploration and formalization. They also raise a more fundamental question. As mentioned above, marginal likelihoods incorporate a component corresponding to the total cost of fitting, or equivalently, to the total learning work done upon going from the prior to the posterior. The use of hierarchical priors in Bayesian model averaging, by borrowing strength across replicated items, essentially reduces the distance between the prior and the posterior at the level of the replicated items, and thus reduces the cost of fitting. However, it is not clear whether it suppresses this cost entirely. If not, then this suggests that Bayesian model averaging might have a general tendency to be over-penalizing, compared to what could be achieved using more aggressive non-Bayesian model fitting approaches.

Cross-validation and wAIC: numerical considerations

In contrast to the marginal likelihood, cross-validation appears to be relatively well-behaved, if the aim is to select the most accurate model. However, it requires some care, both for defining the specific details of the CV procedure and for implementing a reliable numerical approach. In this respect, k-fold CV gives reasonable results, but it is impractical. There are numerical issues with the naive importance sampling approach, which can lead to a serious underestimation of the CV score, in particular for higher-dimensional models.

The k-fold CV approach implemented by naive IS has been used previously for comparing site-heterogeneous and site-homogeneous models (e.g. Philippe *et al.*, 2011; Pisani *et al.*, 2015; Simion *et al.*, 2017). In most cases, site-heterogeneous models have been found as the best fitting models. Importantly, the effective bias of nIS is in favor of less parameter-rich models, which suggests that the fit of the site-heterogeneous models has been underestimated thus far.

The alternative numerical approach used here for computing the k-fold CV score, based on

sequential importance sampling, is much more accurate than nIS. However it is computationally prohibitive. Of note, there are more sophisticated approaches than the one recruited here for implementing sIS, based on particle filters (Wang *et al.*, 2016), which have better Monte Carlo properties than the naive version explored here. However, the single-site importance sampling variance observed here suggests that such particle filters will require many particles and will thus be computationally expensive on data sets of realistic size.

Leave-one-out cross-validation stands out as the computationally most efficient and most easily implemented cross-validation approach. LOO-CV also has a clear asymptotic connection with information criteria, and more specifically with wAIC in the Bayesian case. For large datasets, wAIC is easily implemented. It has a broad range of applicability, at least among i.i.d. models, and the control of its numerical error is easier than for LOO-CV. In the end, based on the results presented here, the practical recommendations for finding the best approximating model are relatively simple. LOO-CV and the wAIC represent the most practical and reliable approaches. Both can be obtained from a single pass over the MCMC sample, but LOO-CV may require larger MCMC sample sizes (or less thinning) than wAIC to pass the quality checks. Thus, if the dataset is sufficiently large (> 5000 aligned positions), the wAIC can be used by default. If, on the other hand, the dataset is small, then LOO-CV should be preferred.

Practical consequences and perspectives

The arguments exposed here in favor of LOO-CV and wAIC over Bayes factors for model approximation purposes are at odds with the general perception in the applied Bayesian community that Bayes factors represent a general gold standard for model selection (Kass & Raftery, 1995; Lartillot & Philippe, 2006; Xie *et al.*, 2011; Oaks *et al.*, 2019). This raises the question of the practical consequences of the use of Bayes factors thus far, in situations where cross-validation might have represented a logically more adequate criterion. As illustrated by the analysis of the normal case (Figure 1), for large datasets and for models that don't differ too much in their dimensionality, all model selection approaches agree in their selection. Thus, in practice, previous results based on the application of Bayes factors on large datasets, such as multi-gene phylogenetic analyses, are unlikely to be qualitatively incorrect, although the case is less clear for smaller-scale analyses. In any case, perhaps a more fundamental contribution of the present analysis is just to facilitate Bayesian model selection, by providing simple guidelines, but also, by making the computational problem of accurately estimating marginal likelihoods practically less relevant.

Still, one main limitation of the methods explored here is that they are valid only in the context of i.i.d. models. For more complex model design, in particular with gene-specific effects in a multi-gene analysis (Suchard *et al.*, 2003; Fan *et al.*, 2011), or when combining sequence data, fossil information and phenotypic or life-history traits (Lartillot & Poujol, 2011; Zhang *et al.*, 2016; Gavryushkina *et al.*, 2016), it is less obvious how to define a correct model selection approach and a meaningful asymptotic analysis.

Of note, this limitation also concerns information criteria classically used in a maximum likelihood framework. There have been many illegitimate applications of criteria such as AIC (or BIC) in non i.i.d. settings. Modified versions of these two criteria have been proposed for the specific case of partition models (Seo & Thorne, 2018; Susko & Roger, 2020). It would be useful to develop a wAIC equivalent of the modified AIC that was proposed in this context and, more generally, in the context of other non-iid settings.

Finally, and apart from the practical considerations, the asymptotic theory behind the development of information criteria such as wAIC also suggests an interesting frequentist perspective on Bayesian inference, opening to more general questions, such as efficiently estimating the sampling bias, variance, error, and more general measures of the frequentist risk of the Bayesian estimators, all of which are worth further exploration.

Materials & Methods

Definitions and relations between alternative Bayesian measures of model fit

In this subsection, the alternative Bayesian measures of model fit are formally defined. A homogeneous mathematical notation is introduced, so as to emphasize the connections and the differences between them, leaving aside in a first step the numerical and algorithmic problems.

Suppose we have a dataset made of n observations, $X = (X_i)_{i=1..n}$. In the context of phylogenetic inference, these observations would typically be the columns of a multiple sequence alignment. In the following, we will adopt a frequentist perspective and assume that these observations are iid from an infinite population of unknown distribution.

We then consider a model M parameterized by θ . In a Bayesian framework, this model is endowed with a prior $p(\theta)$ and then conditioned on data X , giving the posterior distribution

$p(\theta | X)$:

$$p(\theta | X) = \frac{p(X | \theta)p(\theta)}{p(X)} \quad (9)$$

where

$$p(X) = \int p(X | \theta)p(\theta)d\theta \quad (10)$$

is the marginal likelihood.

We wish to evaluate the fit of the model. A first approach is to use the marginal likelihood as the measure of the fit. In the following, when multiple models are compared, the dependence of the marginal likelihood on the specific model will be more explicitly noted $p(X | M)$, for model M . Otherwise, the simpler notation $p(X)$ is used. Often, the fit of a given model M_2 is computed relatively to another model M_1 , by computing the Bayes factor, defined as the ratio of the marginal likelihoods of the two models (Jeffreys, 1935):

$$\text{BF} = \frac{p(X | M_2)}{p(X | M_1)}. \quad (11)$$

A Bayes factor greater than 1 thus means that model M_2 has a higher fit, compared to M_1 . As a way to ensure a scaling consistency across all alternative measures of fit considered here, it is useful to define the per-site log marginal likelihood:

$$m = \frac{1}{n} \ln p(X) \quad (12)$$

or, when comparing two models, the per-site log Bayes factor:

$$\Delta m = \frac{1}{n} (\ln p(X | M_2) - \ln p(X | M_1)) \quad (13)$$

An alternative to marginal likelihoods and Bayes factors is cross-validation. As mentioned in the introduction, the general idea is to split the dataset into two subsets, using one subset (noted X^t) for training the model and then evaluating the fit of the model over the remaining subset (noted X^v , for validation). In the context of Bayesian inference, a natural procedure to implement cross-validation is to average the likelihood under the validation set over the posterior distribution under the training set, i.e. computing:

$$p(X^v | X^t) = \int p(X^v | \theta) p(\theta | X^t) d\theta. \quad (14)$$

The resulting cross-validation score is then log-transformed and averaged over multiple random splits of the original dataset into a training and a validation sets. Of note, alternative approaches

have been proposed, such as computing the cross-validated likelihood on a plug-in estimate, typically, the posterior mean (Konishi & Kitagawa, 1996). However, this approach is not invariant by reparameterization. It is also not applicable for singular and redundant models, such as mixture models (Plummer, 2008).

Based on this general idea, multiple settings can be contemplated for implementing cross-validation. These alternative settings differ in how the dataset is split, how the replication procedure is defined, or whether the likelihood is averaged over the posterior distribution jointly for all observations of the validation set, or independently for each of them.

In k -fold cross-validation, the dataset is split into k subsets of equal size. Then, each subset is considered in turn as the validation set, while the other $k - 1$ subsets are pooled together to make the training set. There are thus k replicates in total. In a variant of this approach, used in phylogenetics (Lartillot & Philippe, 2008), each replicate is obtained independently of other replicates, by randomly splitting the dataset into a fraction f of the observations, which is set aside for validation, while the remaining fraction $1 - f$ used for training. It is thus close to the original version of k -fold cross-validation, with $f = 1/k$, except that the replicates are not obtained by systematic rotation of the subsets. As a result, the number of replicates can be arbitrary. In practice, for computational reasons, a small number of replicates is used, typically $m = 10$. The fraction f is typically set to 0.1 or 0.2, or equivalently, $k = 10$ or 5. In the following, this approach will also be called k -fold cross-validation, even if it does not exactly correspond to the original version.

To more formally describe cross-validation, assume that $l = 1..L$ replicates are considered, each based on a random split of the dataset into $X = (X_l^t, X_l^v)$, and that the training and validation sets are of size q and r , respectively. Thus, in k -fold CV, $q = (1 - f)n$ and $r = fn$, with $f = 1/k$, but the definitions introduced below are valid for more general settings. Using these notations, the final cross-validation score can be defined as:

$$cv_j = \frac{1}{r} \frac{1}{L} \sum_l \ln p(X_l^v | X_l^t) \quad (15)$$

Note that, in this definition, the logarithmic score is divided by the size of the validation set: it is thus a measure of the predictive score per future observation. This definition will be useful for comparing the alternative settings introduced below, which differ in the value of r .

The definition just given, which corresponds to how cross-validation was implemented previously in a phylogenetic context (Lartillot *et al.*, 2007; Lartillot & Philippe, 2008), averages the *joint*

likelihood of all observations of the validation set on the posterior distribution under the training set. Alternatively, the posterior averaging can be done independently for each observation of the validation set. Noting I_l^v the subset of $1..n$ corresponding to the indices of the observations assigned to the validation set in replicate l , the *sitewise* cross-validation score is thus defined as:

$$cv_s = \frac{1}{r} \frac{1}{L} \sum_l \sum_{i \in I_l^v} \ln p(X_i | X_l^t) \quad (16)$$

Again, the score is per future observation. The sitewise approach has apparently not been used previously in the context of phylogenetics. However, it can be useful to consider it in the general comparative evaluation conducted below.

Finally, in *leave-one-out* cross-validation, each observation is taken in turn and set aside for validation, using the $n - 1$ remaining observations to train the model. Noting $X_{(i)}$ the training set (of size $n - 1$) obtained by removing observation i , the score is defined as:

$$cv_l = \frac{1}{n} \sum_n \ln p(X_i | X_{(i)}) \quad (17)$$

Whichever setting is used, from a frequentist perspective, the primary quantity of interest, implicitly targeted by cross-validation, is the expected log-likelihood cross-validation score, the expectation being taken over multiple independent draws of the entire dataset X from the population. Again assuming that the training set X^t is of size q and the validation set X^v is of size r under the chosen setting, this expected score only depends on q and r and is noted:

$$C(q, r) = \frac{1}{r} E [\ln p(X^t | X^v)] \quad (18)$$

Thus, by this definition, k -fold joint cross-validation can be seen as an estimator of $C((1 - f)n, fn)$, with $f = 1/k$, k -fold site-wise cross-validation as an estimator of $C((1 - f)n, 1)$, and leave-one-out cross-validation as an estimator of $C(n - 1, 1)$. Since $C(0, n) = \frac{1}{n} E[\ln p(X)]$, the per-site log marginal likelihood can also be seen as a special case of this formula.

If the model is regular, then for large n , the posterior distribution becomes increasingly concentrated around an asymptotic parameter value θ_0 . In the specific case where the data have been produced under the model, then θ_0 will be the true parameter value. In the general case where the data are from an unknown distribution, there is no true parameter value, in which case θ_0 is the best approximation (in the Kullback-Leibler metric) that the model can give for the distribution induced by this empirical source. In both cases, for large n , all expected scores introduced above, k -fold, leave-one-out or marginal, converge asymptotically to the expected log likelihood of

an observation sampled from the population under θ_0 :

$$C(0, n) \sim C((1-f)n, fn) \sim C((1-f)n, 1) \sim C(n-1, 1) \sim E[\ln p(X_1 | \theta_0)] \quad (19)$$

This justifies the rescaling conventions introduced above.

Monte Carlo methods

To ease notation, in the following, it is assumed that, for cross validation, the first q data points were used for training, and the last r for validation, with q and r depending on the exact cross-validation approach. Estimation of the marginal likelihood can be seen as a special case obtained by setting $q = 0$. In what follows, $X_{a:b}$ denotes the set of observations $(X_i)_{a \leq i \leq b}$. Thus, in particular, $X_{1:q}$ represents the first q observations (i.e. the training set). When $q = 0$, $X_{1:q}$ is the empty set. The index $i = 1..n$ runs over data points, and $k = 1..K$ over the parameter configurations sampled by MCMC.

Naive importance sampling (nIS) for k-fold CV

The naive importance sampling (nIS) approach is used for joint and sitewise k-fold CV. In both cases, we assume that an MCMC chain has been run under the training set, yielding a sample $(\theta_k)_{k=1..K}$ approximately under the posterior distribution $\theta_k \sim p(\theta | X_{1:q})$, for $k = 1..K$.

First considering joint k-fold cross-validation, equation 14, being an expectation over the posterior distribution under the training set, can be approximated by the corresponding Monte Carlo average:

$$p(X_{q+1:n} | X_{1:q}) \simeq \frac{1}{K} \sum_{k=1}^K p(X_{q+1:n} | \theta_k). \quad (20)$$

Thus, nIS for joint k-fold CV runs as follows: for $k = 1..K$, compute the likelihood of the validation data, $L_k = p(X_{q+1:n} | \theta_k)$, compute the arithmetic mean of the L_k 's over the K Monte Carlo samples and log-transform.

A similar approach can be used for the site-wise version of k-fold CV, since, for any single observation X_i of the validation set:

$$p(X_i | X_{1:q}) \simeq \frac{1}{K} \sum_{k=1}^K p(X_i | \theta_k). \quad (21)$$

The site-wise posterior averages can be computed in parallel for each observation and then combined according to equation 16. That is, for $k = 1..K$, compute the likelihood separately for each data

point of the validation data, $L_{ik} = p(X_i | \theta_k)$, for $i = q + 1..n$. In a second step, for all $i = q + 1..n$, compute the arithmetic mean of the L_{ik} 's over the K Monte Carlo samples, log-transform, and finally, sum all individual contributions across the r data points of the validation set.

The cross-predictive ordinate (CPO) approach for LOO-CV

The cross-predictive-ordinate (CPO) approach (Chen *et al.*, 2012; Lewis *et al.*, 2014) gives an estimate of the leave-one-out cross validation score. It relies on the following harmonic-mean identity:

$$\frac{1}{p(X_i | X_{(i)})} = \int \frac{1}{p(X_i | \theta)} p(\theta | X) d\theta \quad (22)$$

This identity suggests to obtain a sample of parameter configurations from the posterior distribution under the entire dataset, $\theta_k \sim p(\theta | X_{1..n})$, for $k = 1..K$, and then approximate the expectation given by equation 22 by a Monte Carlo average:

$$\frac{1}{p(X_i | X_{(i)})} \simeq \frac{1}{K} \sum_{k=1}^K \frac{1}{p(X_i | \theta_k)} \quad (23)$$

Here also, like for site-wise k-fold CV, the Monte Carlo averages across all sites can be computed in parallel, over a single scan of the MCMC chain. Thus, for each $k = 1..K$; for each $i = 1..n$, compute the likelihood of each data point separately, $L_{ik} = p(X_i | \theta_k)$ for site i . Then, in a second step, for each site, compute the harmonic mean \bar{L}_i of the L_{ik} 's over the K Monte Carlo samples, log-transform and sum all individual contributions across the r data points of the validation set.

Sequential Importance Sampling (sIS) for k-fold CV and BF

Sequential Importance sampling is a step-by-step version of IS, which is based on the observation that the joint probability of the validation set can be expressed in terms of a sequential product of the marginal likelihoods of each of individual observations:

$$p(X^v | X^t) = \prod_{i=q+1}^n p(X_i | X_{1:i-1}) \quad (24)$$

or, on a logarithmic scale and on a per-site basis:

$$\frac{1}{r} \ln p(X^v | X^t) = \frac{1}{r} \sum_{i=q+1}^n \ln p(X_i | X_{1:i-1}) \quad (25)$$

Of note, in the case where $q = 0$ and $r = n$, i.e. when the training set is empty and the validation set is the complete original dataset, then equation 25 gives the logarithm of the marginal likelihood

(per site). In turn, if, for $k = 1..K$, θ_{ik} is sampled from the partial posterior based on the first $i - 1$ data points, i.e. $\theta_{ik} \sim p(\theta | X_{1:i-1})$, then an importance sampling estimate of $p(X_i | X_{1:i-1})$ is given by:

$$p(X_i | X_{1:i-1}) \simeq \frac{1}{K} \sum_{k=1}^K p(X_i | \theta_{ik}) \quad (26)$$

This suggests to run a quasi-static MCMC in which data points are added sequentially, each time running the MCMC for a few cycles and averaging the likelihood of the next data point under parameter configurations sampled from the posterior induced by all current data points (equation 26). These individual IS estimates can then be log-transformed and combined additively (equation 25).

To formalize this, in the following, a cycle is defined as a coordinated series of multiple MCMC moves that are applied successively on all parameter components of the models. A parameter configuration is saved after each cycle. A cycle can be arbitrary, although in practice, for sIS to give accurate estimates, a cycle should be sufficiently long to give a reasonably good de-correlation of the MCMC between successive saved samples. The algorithm then proceeds as follows. Starting from a parameter configuration sampled from the prior $\theta_0 \sim p(\theta)$, at step $i = 1..n$:

- the MCMC is run for a short burn-in period of B cycles, so as to equilibrate the MCMC, and then for another series of K cycles, giving K new parameter configurations θ_{ik} approximately under the partial posterior distribution $p(\theta | X_{1:i-1})$;
- the likelihood of the next data point is calculated under each of these sampled parameter values, i.e. $L_{ik} = p(X_i | \theta_{ik})$
- the arithmetic mean of the K likelihood factors L_{ik} , $k = 1..K$, is calculated:

$$L_i = \frac{1}{K} \sum_{k=1}^K L_{ik} \quad (27)$$

- finally, the L_i 's for $i = 1..n$ are log-transformed and combined such as specified by equations 25 and 26.

The quality of the estimate of $p(X_i | X_{1:i-1})$ given by equation 27 depends on K , the number of samples, but also on the variance of the log-likelihood $\ln p(X_i | \theta)$ under the partial posterior $p(\theta | X_{1:i-1})$. In the following, this variance is noted v_i . When this variance is large, a larger value of K should be used. In practice, many data points (e.g. constant sites in a phylogenetic context)

are characterized by a small variance, while a minority of data points induce a large variance. This suggests that the number of Monte Carlo samples K can be tuned on a per-site basis, by first running a small number of cycles at step i to estimate the variance of the log-likelihood and then proceed with a number of cycles K_i determined based on this variance estimate. A heuristic argument for deciding how K_i should scale as a function of v_i is as follows. If the L_{ik} 's were normally distributed, of variance v_i , then the variance of $\ln L_i$, i.e. the log-transformed Monte Carlo estimate given by equation 27 would scale as $V \sim \frac{e^{v_i}}{K_i}$. Inverting this equation, this suggests that, in order to target a given fixed variance for all data points, K_i should scale as e^{v_i} . Taken together, these observations lead to a second version of the algorithm, which runs as follows. Starting from a parameter configuration sampled from the prior $\theta_0 \sim p(\theta)$, at step $i = 1..n$:

- the MCMC is run for B cycles, giving B new parameter configurations θ_{ik} approximately under the partial posterior distribution $p(\theta | X_{1:i-1})$. For each of these K parameter configurations, the likelihood of the next data point is calculated under this sampled parameter value, i.e. $L_{ik} = p(X_i | \theta_{ik})$, for $k = 1..B$.
- the sample variance \hat{v}_i of the $\ln L_{ik}$'s is computed and used to determine K_i , using the following rule: $K_i = \min(K_0 e^{\hat{v}_i - v_0}, K_{max})$, which thus implements an exponential scaling of K_i as a function of v_i , targeting a sample size of K_0 for sites for which have a variance equal to v_0 , and truncated at K_{max} .
- the MCMC is run for another series of K_i cycles, yielding a new sample of parameter configurations θ'_{ik} , for $k = 1..K_i$; the likelihood factors for the next data point are computed, $L'_{ik} = p(X_i | \theta'_{ik})$, for $k = 1..K_i$.
- this second series of likelihoods is used to compute the arithmetic average L_i

Of note, the preliminary run of B cycles at the beginning of each step also contributes to equilibrating the MCMC just after the addition of the last data point. In some cases, there is still a small minority of data points for which v_i may be too large, such that $K_i = K_{max}$ and the Monte Carlo error may not be well controlled. On the other hand, if they represent a small fraction of the data, their contribution to the total error should be small. This point is checked in a second step, based on the effective sample size (Supplementation Information, section 1).

Implementation of nIS, sIS and CPO under the normal model

In the case of the normal model, it is possible to sample the parameter θ directly from the posterior distribution, $p(\theta | X_{1:i})$ for any i (see Appendix). The Monte Carlo implemented for the normal model takes advantage of this property, by sampling θ_k , $k = 1..K$ (for nIS and CPO) or θ_{ik} , for $i = 1..n$ and $k = 1..K$ (for sIS) independently from the relevant posterior distribution.

Implementation of nIS, sIS and CPO under the phylogenetic models

In the case of phylogenetic models, the implementation of PhyloBayesMPI was taken as a starting point. The basic routines of MCMC sampling defining a cycle were left unchanged. Naive importance sampling was already implemented for joint k-fold CV, as a simple post-analysis routine that scans the MCMC chain (after burnin) and averages the likelihood scores over the run. This routine was augmented to also output the sitewise k-fold CV, according to the method described above. Similarly, LOO-CV and the wAIC are jointly computed based on another post-analysis routine, by scanning the MCMC sample, burn-in excluded, computing and storing the site-wise likelihood scores, and finally computing the harmonic mean separately for each site (for LOO-CV), and the arithmetic mean and variance for each site (for wAIC).

The sIS method requires more specific additions to the current implementation: essentially, defining and implementing the family of reference distributions that are necessary for the variance reduction approach, and implementing the routines for adding / removing sites during the MCMC, computing the log-likelihood of a single site, as well as its sample variance, on the fly in order to determine the number of cycles to use (see supplementary information).

General settings across all experiments

For the experiments shown in Table 1, under the normal model, 100 datasets of size $n = 1000$ observations were simulated under the following parameter values: dimension $p = 100, 300, 1000$, true mean $\theta_* = 0.04$, and variance parameter $\sigma^2 = 1$. For each replicate, the analytical values of the log Bayes factor and the relative cross-validation score between the two models were computed (see Appendix). In the case of cross-validation, since the results are ultimately averaged over the 100 independent simulation replicates, only one choice for the splitting of the dataset into a validation and a training set is considered for each replicate, taking the first 800 data points for training and the last 200 for validation. For sIS, the first approach (fixed K for all data points) was used, with

K such as indicated in Table 1.

For the phylogenetic analyses, two previously published empirical datasets were considered:

- EF2: a multiple sequence alignment of elongation factor 2 in 30 eukaryotic species (627 aligned positions), taken from Lartillot & Philippe (2006);
- Metazoa: and a concatenation of genes (35371 aligned positions) across 132 metazoans, along with 2 choanoflagellates and 12 fungi for the outgroup (Philippe *et al.*, 2005).

In the case of EF2 (Table 2), for k -fold CV (both joint and site-wise), 10 random reshufflings of the sites of the original dataset were created. Each of these reshuffled version of the dataset was further split into a training and a validation set, on which nIS was applied, running the MCMC for 1100 points (saving every cycle for JTT, LG, GTR and every 3 cycles for CAT) on the training set and then computing the CV score of the validation set based on the 1000 last points of the MCMC. For sIS, a preliminary run on the original dataset and under each model was conducted for 1100 points (saving every 3 cycles for CAT). Empirical posterior means and variances for defining the reference priors for sIS were obtained based on the last 1000 points of this preliminary run. The self-tuned version of the sIS method was then used, with the following parameter values: $B = 10$, $K_0 = 30$, $K_{max} = 1000$ and $v_0 = 0.1$, with $L = 2$ independent runs for each data replicate (using the original dataset for computing the marginal likelihood and the 10 random reshufflings of the sites for k -fold CV). For LOO-CV, $L = 2$ independent runs were conducted on the complete dataset and under each model, again for 1100 points and saving every 3 cycles for CAT, discarding the first 100 points of burn-in and applying the CPO method on the 1000 remaining points.

For the experiments shown in Figure 2, the Metazoan dataset was first filtered to remove sites with more than 20% of missing data, leaving a total of 9804 sites. Then, 4 jackknife replicates of sizes $n = 200$ to 800 were randomly sampled. For figure 3, a standard MCMC was run under the complete dataset (9804 sites) and under the LG model, for a total of 1100 points. Then, 10 posterior predictive replicates were simulated, based on 10 samples regularly spaced across this MCMC chain, discarding the first 100 points and taking one every 100 points. Each of these simulated datasets was then jackknifed, yielding replicates of sizes $n = 200$ to $n = 800$. Both the empirical and the simulated jackknife replicates were then used for assessing the fit of the model by marginal likelihood or LOO-CV, using the same settings as for EF2.

For Figure 4, the complete dataset was used, and 4 jackknife replicates of size ranging from $n = 200$ to $n = 16000$ were randomly sampled. The LOO-CV scores and the two versions of the

wAIC were computed by running $L = 2$ independent chains of 1100 points (saving every 3 cycles for CAT) on each jackknife replicate, using the last 1000 points, with or without a 10-fold thinning (in which case the Monte Carlo estimates of the posterior averages are based on 100 points).

Data and software availability

All methods for phylogenetic models were implemented in PhyloBayes (Lartillot *et al.*, 2013), version 1.9, available at <https://github.com/bayesiancook/pbmpi>. The empirical data used here are also available through this repo, along with example scripts. The methods for the analytical and Monte Carlo results under the normal model are available at <https://github.com/bayesiancook/normcv>.

Competing interests

The Author declares no competing interest.

Acknowledgements

The Author wishes to thank Marie-Laure Delignette, Philippe Veber, Nicolas Rodrigue and Hervé Philippe for discussions and thoughtful comments on the manuscript.

Funding

This work was granted access to the HPC resources of CINES under the allocation A0040310449 made by GENCI, and to the computing cluster PRABI-LBBE.

References

- Aho, K., Derryberry, D. & Peterson, T. 2014 Model selection for ecologists: the worldviews of AIC and BIC. *Ecology*, **95**(3), 631–636.
- Akaike, H. 1974 A new look at the statistical model identification. *IEEE Trans. Automat. Contr.*, **19**(6), 716–723.
- Baele, G. & Lemey, P. 2013 Bayesian evolutionary model testing in the phylogenomics era: matching model complexity with computational efficiency. *Bioinformatics*, **29**(16), 1970–1979.

- Baele, G., Lemey, P., Bedford, T., Rambaut, A., Suchard, M. A. & Alekseyenko, A. V. 2012*a* Improving the accuracy of demographic and molecular clock model comparison while accommodating phylogenetic uncertainty. *Mol. Biol. Evol.*, **29**(9), 2157–2167.
- Baele, G., Lemey, P. & Vansteelandt, S. 2013 Make the most of your samples: Bayes factor estimators for high-dimensional models of sequence evolution. *BMC Bioinformatics*, **14**, 85.
- Baele, G., Li, W. L. S., Drummond, A. J., Suchard, M. A. & Lemey, P. 2012*b* Accurate Model Selection of Relaxed Molecular Clocks in Bayesian Phylogenetics. *Mol. Biol. Evol.*
- Bartlett, M. S. 1957 A comment on D. V. Lindley’s statistical paradox. *Biometrika*, **44**, 533–534.
- Berger, J. 2006 The case for objective Bayesian analysis. *Bayesian Analysis*, **1**(3), 385–402.
- Berger, J. O. 1985 *Statistical Decision Theory and Bayesian Analysis*. New-York: Springer-Verlag, 1985th edn.
- Bernardo, J. M. & Smith, A. F. M. 1994 *Bayesian theory*. Chichester, UK: John Wiley & Sons, Inc.
- Blanquart, S. & Lartillot, N. 2006 A Bayesian compound stochastic process for modeling nonstationary and nonhomogeneous sequence evolution. *Mol. Biol. Evol.*, **23**(11), 2058–2071.
- Breiman, L., Friedman, J., Stone, C. J. & Olshen, R. A. 1984 *Classification and Regression Trees*. Taylor & Francis.
- Brown, J. M. & Thomson, R. C. 2017 Bayes Factors Unmask Highly Variable Information Content, Bias, and Extreme Influence in Phylogenomic Analyses. *Syst. Biol.*, **66**(4), 517–530.
- Burnham, K. P. & Anderson, D. R. 2002 *Model Selection and Multimodel Inference: a practical information-theoretic approach*. New-York: Springer, 2nd edn.
- Celeux, G., Forbes, F., Robert, C. P. & Titterton, D. M. 2006 Deviance Information Criteria for Missing Data Models. *Bayesian Analysis*, **1**(4), 651–674.
- Chen, M. H., Shao, Q. M. & Ibrahim, J. G. 2012 *Monte Carlo Methods in Bayesian Computation*. Springer Series in Statistics. Springer New York.
- Efron, B. 1986 How Biased is the Apparent Error Rate of a Prediction Rule? *Journal of the American Statistical Association*, **81**(394), 461–470.
- Evans, J. & Sullivan, J. 2012 Generalized mixture models for molecular phylogenetic estimation. *Syst. Biol.*, **61**(1), 12–21.

- Fan, Y., Wu, R., Chen, M.-H., Kuo, L. & Lewis, P. O. 2011 Choosing among partition models in Bayesian phylogenetics. *Mol. Biol. Evol.*, **28**(1), 523–532.
- Fragoso, T. M., Bertoli, W. & Neto, F. L. 2017 Bayesian model averaging: A systematic review and conceptual. *International Statistical Review*, **86**(1), 1–28.
- Gavryushkina, A., Heath, T. A., Ksepka, D. T., Stadler, T., Welch, D. & Drummond, A. J. 2016 Bayesian Total-Evidence Dating Reveals the Recent Crown Radiation of Penguins. *Syst. Biol.*
- Geisser, S. 1975 The predictive sample reuse method with application. *Journal of the American Statistical Association*, **70**, 320–328.
- Geisser, S. & Eddy, W. 1979 A predictive approach to model selection. *Journal of the American Statistical Association*, **74**, 153–160.
- Gelfand, A. E. 1996 Model determination using sampling-based methods. In *Markov chain monte carlo in practice* (eds W. R. Gilks, S. Richardson & D. J. Spiegelhalter), pp. 145–162. Chapman & Hall/CRC.
- Gelfand, A. E., Dey, D. K. & Chang, H. 1992 Model determination using predictive distributions with implementation via sampling-based methods. In *Bayesian statistic, 4th edn* (eds J. M. Bernardo, J. O. Berger, A. P. Dawid & A. F. M. Smith), pp. 147–167. Oxford: Oxford University Press.
- Gelman, A., Hwang, J. & Vehtari, A. 2014 Understanding predictive information criteria for Bayesian models. *Stat Comput*, **24**, 997–1016.
- Goldman, N. 1993 Statistical tests of models of DNA substitution. *J. Mol. Evol.*, **36**(2), 182–198.
- Hoeting, J. A., Madigan, D., Raftery, A. E. & Volinsky, C. T. 2000 Bayesian Model Averaging: A Tutorial. *Statistical Science*, **14**(4), 382–417.
- Huelsenbeck, J. P., Larget, B. & Alfaro, M. E. 2004 Bayesian phylogenetic model selection using reversible jump Markov chain Monte Carlo. *Mol. Biol. Evol.*, **21**(6), 1123–1133.
- Jeffreys, H. 1935 Some tests of significance, treated by the theory of probability. *Proc. Camb. Phil. Soc.*, **31**, 203–222.
- Jeffreys, H. 1967 *Theory of probability*. London: Oxford University Press.
- Jones, D. T., Taylor, W. R. & Thornton, J. M. 1992 The rapid generation of mutation data matrices from protein sequences. *Comput Appl Biosci*, **8**(3), 275–282.
- Kass, R. E. & Raftery, A. E. 1995 Bayes Factors. *Journal of the American Statistical Association*, **90**(430), 773–795.

- Konishi, S. & Kitagawa, G. 1996 Generalised information criteria in model selection. *Biometrika*, **83**(4), 875–890.
- Konishi, S. & Kitagawa, G. 2007 *Information Criteria and Statistical Modeling*. Springer New York.
- Kosakovsky Pond, S. L. & Frost, S. D. W. 2005 Not so different after all: a comparison of methods for detecting amino acid sites under selection. *Mol. Biol. Evol.*, **22**(5), 1208–1222.
- Koshi, J. M. & Goldstein, R. A. 2001 Analyzing site heterogeneity during protein evolution. *Pac Symp Biocomput*, pp. 191–202.
- Lartillot, N., Brinkmann, H. & Philippe, H. 2007 Suppression of long-branch attraction artefacts in the animal phylogeny using a site-heterogeneous model. *BMC Evol. Biol.*, **7 Suppl 1**, S4.
- Lartillot, N., Lepage, T. & Blanquart, S. 2009 PhyloBayes 3: a Bayesian software package for phylogenetic reconstruction and molecular dating. *Bioinformatics*, **25**(17), 2286–2288.
- Lartillot, N. & Philippe, H. 2004 A Bayesian mixture model for across-site heterogeneities in the amino-acid replacement process. *Mol. Biol. Evol.*, **21**(6), 1095–1109.
- Lartillot, N. & Philippe, H. 2006 Computing Bayes factors using thermodynamic integration. *Syst. Biol.*, **55**(2), 195–207.
- Lartillot, N. & Philippe, H. 2008 Improvement of molecular phylogenetic inference and the phylogeny of Bilateria. *Philos. Trans. R. Soc. Lond. B Biol. Sci.*, **363**(1496), 1463–1472.
- Lartillot, N. & Poujol, R. 2011 A phylogenetic model for investigating correlated evolution of substitution rates and continuous phenotypic characters. *Mol. Biol. Evol.*, **28**(1), 729–744.
- Lartillot, N., Rodrigue, N., Stubbs, D. & Richer, J. 2013 PhyloBayes MPI: phylogenetic reconstruction with infinite mixtures of profiles in a parallel environment. *Syst. Biol.*, **62**(4), 611–615.
- Le, S. Q. & Gascuel, O. 2008 An improved general amino acid replacement matrix. *Mol. Biol. Evol.*, **25**(7), 1307–1320.
- Lewis, P. O., Xie, W., Chen, M.-H., Fan, Y. & Kuo, L. 2014 Posterior predictive Bayesian phylogenetic model selection. *Syst. Biol.*, **63**(3), 309–321.
- Lindley, D. V. 1957 A statistical paradox. *Biometrika*, **44**, 187–192.
- Nielsen, R. & Yang, Z. 1998 Likelihood models for detecting positively selected amino acid sites and applications to the HIV-1 envelope gene. *Genetics*.

- Oaks, J. R., Cobb, K. A., Minin, V. A. & Leaché, A. D. 2019 Marginal Likelihoods in Phylogenetics: A Review of Methods and Applications. *Syst. Biol.*, **68**(5), 681–697.
- Pagel, M. & Meade, A. 2004 A phylogenetic mixture model for detecting pattern-heterogeneity in gene sequence or character-state data. *Syst. Biol.*, **53**(4), 571–581.
- Philippe, H., Brinkmann, H., Copley, R. R., Moroz, L. L., Nakano, H., Poustka, A. J., Wallberg, A., Peterson, K. J. & Telford, M. J. 2011 Acoelomorph flatworms are deuterostomes related to Xenoturbella. *Nature*, **470**(7333), 255–258.
- Philippe, H., Lartillot, N. & Brinkmann, H. 2005 Multigene analyses of bilaterian animals corroborate the monophyly of Ecdysozoa, Lophotrochozoa, and Protostomia. *Mol. Biol. Evol.*, **22**(5), 1246–1253.
- Pisani, D., Pett, W., Dohrmann, M., Feuda, R., Rota-Stabelli, O., Philippe, H., Lartillot, N. & Wörheide, G. 2015 Genomic data do not support comb jellies as the sister group to all other animals. *Proceedings of the National Academy of Sciences*, **112**(50), 15 402–15 407.
- Plummer, M. 2008 Penalized loss functions for Bayesian model comparison. *Biostatistics*, **9**(3), 523–539.
- Quang, L. S., Gascuel, O. & Lartillot, N. 2008 Empirical profile mixture models for phylogenetic reconstruction. *Bioinformatics*, **24**(20), 2317–2323.
- Raftery, A. E., Newton, M. A., Satagopan, J. M. & Krivitsky, P. N. 2007 Estimating the Integrated Likelihood via Posterior Simulation Using the Harmonic Mean Identity. *Bayesian Statistics*, **8**, 1–45.
- Ronquist, F., Kudlicka, J., Senderov, V., Borgström, J., Lartillot, N., Lundén, D., Murray, L., Schön, T. B. & Broman, D. 2021 Universal probabilistic programming offers a powerful approach to statistical phylogenetics. *Communications Biology*, pp. 1–10.
- Schrempf, D., Lartillot, N. & Szöllösi, G. 2020 Scalable empirical mixture models that account for across-site compositional heterogeneity. *Mol. Biol. Evol.*
- Schwarz, G. 2006 Estimating the Dimension of a Model. *Ann. Statist.*, **6**(2), 461–464.
- Seo, T.-K. & Thorne, J. L. 2018 Information Criteria for Comparing Partition Schemes. *Syst. Biol.*, **67**(4), 616–632.
- Shao, J. 1993 Linear Model Selection by Cross-Validation. *Journal of the American Statistical Association*, **88**(422), 486–494.
- Shibata, R. 1986 Consistency of Model Selection and Parameter Estimation. *Journal of applied probability*, **23**, 127–141.

- Shibata, R. 1989 Statistical aspects of model selection. In *From data to model* (ed. J. C. Willems), pp. 215–240. Springer New York.
- Shimodaira, H. 2004 Approximately unbiased tests of regions using multistep-multiscale bootstrap resampling. *Ann. Statist.*, **32**(6), 2616–2641.
- Simion, P., Philippe, H., Baurain, D., Jager, M., Richter, D. J., Di Franco, A., Roure, B., Satoh, N., Quéinnec, E. *et al.* 2017 A Large and Consistent Phylogenomic Dataset Supports Sponges as the Sister Group to All Other Animals. *Curr. Biol.*, **27**(7), 958–967.
- Smyth, P. 2000 Model selection for probabilistic clustering using cross-validated likelihood - Springer. *Stat Comput.*
- Spiegelhalter, D. J., Best, N. G. & Carlin, B. P. 2002 Bayesian measures of model complexity and fit. *J. R. Statist. Soc. B*, **64**, 583–639.
- Spiegelhalter, D. J., Best, N. G., Carlin, B. P. & Van der Linde, A. 2014 The deviance information criterion: 12 years on. *J. R. Statist. Soc. B*, **76**(3), 485–493.
- Stone, M. 1974 Cross-validatory choice and assessment of statistical predictions. *J. R. Statist. Soc. B*, pp. 111–147.
- Stone, M. 1977 An asymptotic equivalence of choice of model by cross-validation and Akaike’s criterion. *J. R. Statist. Soc. B*, pp. 44–47.
- Suchard, M. A., Kitchen, C. M. R., Sinsheimer, J. S. & Weiss, R. E. 2003 Hierarchical Phylogenetic Models for Analyzing Multipartite Sequence Data. *Syst. Biol.*, **52**(5), 649–664.
- Suchard, M. A., Weiss, R. E. & Sinsheimer, J. S. 2001 Bayesian selection of continuous-time Markov chain evolutionary models. *Mol. Biol. Evol.*, **18**(6), 1001–1013.
- Sullivan, J. & Joyce, P. 2005 Model selection in phylogenetics. *Annu. Rev. Ecol. Evol. Syst.*, pp. 445–466.
- Susko, E., Lincker, L. & Roger, A. J. 2018 Accelerated Estimation of Frequency Classes in Site-Heterogeneous Profile Mixture Models. *Mol. Biol. Evol.*, **35**(5), 1266–1283.
- Susko, E. & Roger, A. J. 2020 On the Use of Information Criteria for Model Selection in Phylogenetics. *Mol. Biol. Evol.*, **37**(2), 549–562.
- Thomas, V., Pedregosa, F., van Merriënboer, B., Mangazol, P.-A., Bengio, Y. & Le Roux, N. 2020 On the interplay between noise and curvature and its effect on optimization and generalization. In *Proceedings of the 23rd international conference on artificial intelligence and statistics (aistats)*.

- Vehtari, A., Gelman, A. & Gabry, J. 2016 Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC. *Stat Comput*, **27**(5), 1413–1432.
- Vrieze, S. I. 2012 Model selection and psychological theory: a discussion of the differences between the Akaike information criterion (AIC) and the Bayesian information criterion (BIC). *Psychol Methods*, **17**(2), 228–243.
- Wang, H.-C., Li, K., Susko, E. & Roger, A. J. 2008 A class frequency mixture model that adjusts for site-specific amino acid frequencies and improves inference of protein phylogeny. *BMC Evol. Biol.*, **8**, 331.
- Wang, L., Bouchard-Côté, A. & Doucet, A. 2016 Bayesian Phylogenetic Inference Using a Combinatorial Sequential Monte Carlo Method. *Journal of the American Statistical Association*, **110**(512), 1362–1374.
- Watanabe, S. 2001 Algebraic geometrical methods for hierarchical learning machines. *Neural Netw*, **14**(8), 1049–1060.
- Watanabe, S. 2007 Almost All Learning Machines are Singular. *IEEE Symposium on Foundations of Computational Intelligence*.
- Watanabe, S. 2009 *Algebraic Geometry and Statistical Learning Theory*. Cambridge Monographs on Applied and Computational Mathematics. Cambridge University Press.
- Watanabe, S. 2010 Asymptotic Equivalence of Bayes Cross Validation and Widely Applicable Information Criterion in Singular Learning Theory. *The Journal of Machine Learning Research*, **11**, 3571–3594.
- Whelan, S. & Goldman, N. 2001 A general empirical model of protein evolution derived from multiple protein families using a maximum-likelihood approach. *Mol. Biol. Evol.*, **18**(5), 691–699.
- Xie, W., Lewis, P. O., Fan, Y., Kuo, L. & Chen, M.-H. 2011 Improving marginal likelihood estimation for Bayesian phylogenetic model selection. *Syst. Biol.*, **60**(2), 150–160.
- Zhang, C., Stadler, T., Klopstein, S., Heath, T. A. & Ronquist, F. 2016 Total-Evidence Dating under the Fossilized Birth-Death Process. *Syst. Biol.*, **65**(2), 228–249.
- Zhang, J., Nielsen, R. & Yang, Z. 2005 Evaluation of an improved branch-site likelihood method for detecting positive selection at the molecular level. *Mol. Biol. Evol.*
- Zhang, P. 1993 Model selection via multifold cross validation. *The Annals of Statistics*, **21**(1), 299–313.