



# Statistical analysis of diversification with species traits

Emmanuel Paradis

## ► To cite this version:

Emmanuel Paradis. Statistical analysis of diversification with species traits. *Evolution - International Journal of Organic Evolution*, 2005, 59 (1), pp.1 - 12. 10.1111/j.0014-3820.2005.tb00889.x . hal-01822136

**HAL Id: hal-01822136**

**<https://hal.science/hal-01822136>**

Submitted on 23 Jun 2018

**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.

STATISTICAL ANALYSIS OF DIVERSIFICATION WITH  
SPECIES TRAITS

EMMANUEL PARADIS

*Laboratoire de Paléontologie, Paléobiologie & Phylogénie,*

*Institut des Sciences de l'Évolution,*

*Université Montpellier II,*

*F-34095 Montpellier cédex 05, France*

*E-mail: paradis@isem.univ-montp2.fr*

LRH: EMMANUEL PARADIS

RRH: STATISTICAL ANALYSIS OF DIVERSIFICATION

*Abstract.*—Testing whether some species traits have a significant effect on diversification rates is central in the assessment of macroevolutionary theories. However, we still lack a powerful method to tackle this objective. I present a new method for the statistical analysis of diversification with species traits. The required data are observations of the traits on recent species, the phylogenetic tree of these species, and reconstructions of ancestral values of the traits. Several traits, either continuous or discrete, and in some cases their interactions, can be analyzed simultaneously. The parameters are estimated by the method of maximum likelihood. The statistical significance of the effects in a model can be tested with likelihood ratio tests. A simulation study showed that past random extinction events do not affect the type I error rate of the tests, whereas statistical power is decreased, though some power is still kept if the effect of the simulated trait on speciation is strong. The use of the method is illustrated by the analysis of published data on Primates. The analysis of these data showed that the apparent overall positive relationship between body mass and species diversity is actually an artefact due to a clade-specific effect. Within each clade the effect of body mass on speciation rate was in fact negative. The present method allows to take both effects (clade and body mass) into account simultaneously.

*Keywords.*—diversification, extinction, maximum likelihood, phylogeny, generalized linear models, speciation.

The connection between models and reality is the central issue.

Another theorem cannot help, because we cannot prove by pure mathematics that mathematical models apply to real problems.

Freedman (1995/96)

Elucidating the mechanisms behind the variation in species diversity is one of the fundamental questions of evolutionary biology. The great disparity in numbers of species among higher taxonomic units (phyla, classes, or orders) implies that there must be some explanatory reasons for this variation. Remarkable progress has been accomplished recently by theoreticians to untangle the variables that affect speciation and/or extinction rates (see Gavrillets 2003, for an overview of some results). Characterizing these variables in empirical studies is not straightforward. Ideally, one would collect data through time using the geological and fossil records, then analyze them with standard statistical methods. However, paleontological and paleoenvironmental data have been found difficult to interpret in this respect (Nichols et al. 1986; Marshall 1997; Bleiweiss 1998; Tipper 1998; Weiss and Marshall 1999; Benton et al. 2000; Foote and Sepkoski 1999; Foote et al. 1999; Jablonski 2000; Archibald and Deutschman 2001).

During the last decade, there has been an increased interest in the use of phylogenies reconstructed from recent species to study evolutionary processes (Barracough and Nee 2001). The characteristics of phylogenetic trees (topology, branch lengths, balance, ...) are affected by the processes of diversification of the clade under consideration (Aldous 1995, 2001; Mooers and Heard 1997). Several methods have been developed to use phylogenetic data in order to estimate speciation and extinction rates (Nee et al. 1994; Paradis 2003), or test for constancy of diversification through time (Wollenberg et al. 1996; Paradis 1997, 1998b; Pybus and Harvey 2000) or among clades (Sanderson and

Bharathan 1993; Sanderson and Donoghue 1994; Paradis 1998a; Bokma 2003).

Other methods aim to test whether some traits are significantly correlated with species diversification (Slowinski and Guyer 1993; Barraclough et al. 1998). This category of methods is of particular relevance since they allow test of hypotheses about the effects of some biological traits on species diversification. Examples of such traits include suspected evolutionary key-innovations (Hunter 1998), body size (Gittleman and Purvis 1998; Bokma 2002; Orme et al. 2002a), or other traits (e.g., Cardillo et al. 2003; Fisher et al. 2003). However, these methods suffer from the fact that they consider the product of evolution (namely species richness) rather than the process of evolution *per se* (speciation and/or extinction). Even analyzing species richness in a phylogenetic framework is not straightforward because of the possibility to use various indices of species richness differences, and the difficulty in combining several traits in the analysis (Isaac et al. 2003).

In this paper, I present a new method to analyze whether some traits affect diversification rate using observations of these traits on recent species and the phylogenetic tree of these species. This method is based on a new model of diversification which can be called the ‘Yule model with covariates’. The parameters of the model are estimated by maximum likelihood. Several traits, either continuous or discrete, can be analyzed simultaneously. The statistical significance of the effects of these traits can be tested with likelihood ratio tests. Since the present method is parametric, it is possible to predict how speciation rate varies with respect to the variables that have been found to be significant.

I develop below the model, the estimation procedure, and the statistical tests. I then report the results from a simulation study to assess the statistical properties of the method. To illustrate the use of the method, I analyzed some data on Primates.

## METHODS

*The Model*

For simplicity, I will assume that the lineage diversify following a birth-only process (i.e. no extinction). In this model, each species has an instantaneous probability of splitting in two daughter-species, called the speciation rate in this paper, and denoted  $\lambda$ . Let us assume that the value of  $\lambda$  for species  $i$  ( $\lambda_i$ ) is determined by a set of species characters. For convenience, I will assume a linear relationship that can be written as:

$$\lambda_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \alpha, \quad (1)$$

where  $x_1, \dots, x_p$  are the species characters under consideration ( $x_{i1}, \dots, x_{ip}$  are the values of these characters for species  $i$ ), and  $\beta_1, \dots, \beta_p, \alpha$  are coefficients. Since  $\lambda_i$  is a probability, it can take values between 0 and 1; it is thus preferable to transform the left-hand side term above so that it varies between  $-\infty$  and  $+\infty$ . Various functions can be used to this aim: I chose a logit function which is widely used in statistical analyses (e.g. in logistic regression):

$$\ln \frac{\lambda_i}{1 - \lambda_i} = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \alpha. \quad (2)$$

It is convenient to rewrite this equation in matrix form:

$$\ln \frac{\lambda_i}{1 - \lambda_i} = x_i^T \beta, \quad (3)$$

where  $\beta$  is a vector composed with the  $p + 1$  coefficients, and  $x_i^T$  is the transposed vector of the  $p$  characters for species  $i$  binded with a 1 for the intercept ( $\alpha$ ). Thus, giving some values of  $x$  and  $\beta$ , the speciation rate is given by the inverse logit function:

$$\lambda_i = \frac{1}{1 + e^{-x_i^T \beta}}. \quad (4)$$

Note that there is no direct time dependence in the speciation rate in this model: if the species characters vary through time, then  $\lambda$  will also vary, but if the former do not, then the latter will be constant through time.

The species characters (or variables) can be continuous or discrete. In the latter case, a suitable numeric coding may be used as common practice in generalized linear modeling (McCullagh and Nelder 1989). For instance, a discrete variable with two states would be coded as a variable taking the values 0 or 1 (see Appendix). Interactions between discrete variables, and between a continuous one and a discrete one can also be included in the model.

### *Parameter Estimation*

The problem is to estimate the vector of coefficients  $\beta$ . I will assume that the available data are a phylogenetic tree that includes all species and with branch lengths proportional to time, and the values of the  $p$  species characters for the  $N$  species and for the  $N - 1$  nodes of the tree (assuming that the latter is fully dichotomous). A series of speciation events can be inferred from the tree; it is also inferred that no speciation occurred between two such consecutive events. Let us denote  $t_1, \dots, t_q$  the times defined by the speciation events, and  $n_1, \dots, n_q$  the number of species living after each of these events ( $n_q = N$ ). The first speciation event is given by the root of the tree, so that  $t_1$  is the age of the root, and  $n_1 = 2$ . I will admit that there may be ties in these times (i.e. speciation events occurring at the same time), so the number of species may increase by more than 1 through time, and  $q$  may be less than  $N - 1$ . The probabilities of these events are proportional to the number of species living at a given time and the corresponding value of speciation rates; if two speciation

events occurred at the same time (i.e. a tie), then this is the product of the number of species and the corresponding values of the two speciation rates (and so on with three or more speciation events).

The probability of a species not having a speciation event between times  $t'$  and  $t''$  is (Cox and Oakes 1984):

$$\exp \left\{ - \int_{t'}^{t''} \lambda(\tau) d\tau \right\}, \quad (5)$$

where  $\lambda(\tau)$  is a function describing how the speciation rate changes between times  $t'$  and  $t''$ . The probability of  $n$  species not having any speciation event between times  $t'$  and  $t''$  is (see Darwin 1956):

$$\exp \left\{ - \sum_{i=1}^n \int_{t'}^{t''} \lambda_i(\tau) d\tau \right\}. \quad (6)$$

We can now write the likelihood of the data under the model described above:

$$\begin{aligned} L = & \exp \left[ - \left\{ \int_{t_1}^{t_2} \lambda_1(\tau) d\tau + \int_{t_1}^{t_2} \lambda_2(\tau) d\tau \right\} \right] n_2 \nu_2 \times \\ & \exp \left[ - \left\{ \int_{t_2}^{t_3} \lambda_1(\tau) d\tau + \int_{t_2}^{t_3} \lambda_2(\tau) d\tau + \int_{t_2}^{t_3} \lambda_3(\tau) d\tau \right\} \right] n_3 \nu_3 \times \\ & \dots \times \\ & \exp \left[ - \left\{ \int_{t_{q-1}}^{t_q} \lambda_1(\tau) d\tau + \dots + \int_{t_{q-1}}^{t_q} \lambda_N(\tau) d\tau \right\} \right], \end{aligned}$$

where  $\nu_2, \dots, \nu_q$  are the probabilities of the events observed at times  $t_2, \dots, t_q$  (namely one or more speciation events). A logarithmic transformation gives:

$$\begin{aligned} \ln L = & - \left\{ \int_{t_1}^{t_2} \lambda_1(\tau) d\tau + \int_{t_1}^{t_2} \lambda_2(\tau) d\tau \right\} + \ln n_2 + \ln \nu_2 + \\ & - \left\{ \int_{t_2}^{t_3} \lambda_1(\tau) d\tau + \int_{t_2}^{t_3} \lambda_2(\tau) d\tau + \int_{t_2}^{t_3} \lambda_3(\tau) d\tau \right\} + \ln n_3 + \ln \nu_3 + \\ & \dots + \\ & - \left\{ \int_{t_{q-1}}^{t_q} \lambda_1(\tau) d\tau + \dots + \int_{t_{q-1}}^{t_q} \lambda_N(\tau) d\tau \right\}. \end{aligned}$$



Using the additive property of integrals, it is easy to see that all integrals add up over all branches. The log-likelihood thus becomes:

$$\ln L = - \sum_{j=1}^{2N-2} \int_{t'_j}^{t''_j} \lambda(\tau) d\tau + \sum_{i=1}^q \ln n_i + \sum_{i=1}^q \ln \nu_i, \quad (7)$$

where  $t'_j$  and  $t''_j$  are the end-point times of the  $j$ th branch ( $2N - 2$  is the number of branch in a rooted dichotomous tree with  $N$  tips). A further simplification is given by the fact that  $\nu$  is a product of  $\lambda$  values, so we have:

$$\ln L = - \sum_{j=1}^{2N-2} \int_{t'_j}^{t''_j} \lambda(\tau) d\tau + \sum_{i=1}^q \ln n_i + \sum_{k=1}^{N-1} \ln \lambda_k. \quad (8)$$

The way  $\lambda(\tau)$  is integrated over the branch lengths will depend on how the speciation rate varies through time. The simplest way is to assume a linear change:

$$\int_{t'_j}^{t''_j} \lambda(\tau) d\tau = \frac{\lambda_{t'_j} + \lambda_{t''_j}}{2} l_j, \quad (9)$$

where  $l_j$  is the length of the  $j$ th branch. Under this assumption, we end up with the log-likelihood:

$$\ln L = - \sum_{j=1}^{2N-2} \frac{\lambda_{t'_j} + \lambda_{t''_j}}{2} l_j + \sum_{i=1}^q \ln n_i + \sum_{k=1}^{N-1} \ln \lambda_k, \quad (10)$$

which can be solved after substituting  $\lambda$  by the expression in equation 4.

Practically, maximizing  $\ln L$  would require to find its first partial derivatives with respect to the different parameters; however these expression are too complex to be solved analytically. I used instead a numerical minimization method. The likelihood was transformed as the deviance (equal to  $-2 \ln L$ ); finding the maximum likelihood was equivalent to minimizing the deviance. This was done with the nonlinear minimization function of R (Ihaka and

Gentleman 1996) which uses a Newton-type algorithm for unconstrained minimization (Schnabel et al. 1985). This method allows one to find the minimum of a function even when its partial and second derivatives are unknown. However, these derivatives can be computed numerically by this algorithm, allowing to calculate the standard-errors of the MLEs with:

$$\text{SE}(\hat{\beta}) = \left( \left[ -\frac{\partial^2 \ln L}{\partial \beta^2} \right]_{\hat{\beta}} \right)^{-\frac{1}{2}}. \quad (11)$$

It is interesting to note that if we fix  $\beta = 0$ , then the estimated speciation rate  $\hat{\lambda}$ , after logit back-transformation of the estimated parameter  $\hat{\alpha}$  using equation 4, is equal to the Kendall–Moran estimate (see Nee 2001, for details), as can be found with the function `yule` in APE (Paradis et al. 2004). The maximum likelihood found by both fitting procedures is the same.

### *Hypothesis Testing*

To test whether the species variables  $x$  have a significant effect on speciation rate, it is possible to use standard statistical methodologies. Two models can be compared with a likelihood ratio test provided they are nested (i.e. one model is a particular case of the other). The hypothesis then tested is that the additional parameter(s) in the second model is (are) significantly different from zero. This likelihood ratio test follows a  $\chi^2$  distribution with a number of degrees of freedom equal to the difference in numbers of parameters between both models.

### *Diagnostics of Fit*

Whereas statistical tests, such as likelihood ratio tests, indicate whether a model better describes the data than another one, they do not inform us whether this particular model adequately fits the data in some absolute ways. Diagnostics of fit are statistical tools to examine whether the observed data are well predicted

by a particular model. For instance in a linear regression, residuals (the difference between the observed and the predicted values of the response) may indicate whether some potential predictors have not been included in the model, or whether the relationship may be nonlinear (Venables and Ripley 2002).

With the method presented in this paper, it is possible to assess whether the branch lengths observed in a tree are well predicted by the Yule model with covariates using equation 5. This equation gives the probability of no speciation event, and so corresponds to the probability of observing a particular branch in the reconstructed tree. Under the above assumption on character change (i.e. linear change), equation 9 can be used practically.

### *Simulation Analysis*

Some simulations were run to assess some statistical properties of the present method. Specifically, four questions were addressed. What is the type I error rate of the method? What is the statistical power of the method? What is the precision of the parameter estimators? What is the statistical robustness of the method with respect to extinction?

The evolution of a clade was numerically simulated using a non-homogeneous birth–death process where the speciation rate was determined by

$$\lambda = \frac{1}{1 + e^{\beta x - \alpha}},$$

where  $\alpha$  and  $\beta$  were parameters, and  $x$  was a continuous variable which evolved following a random walk process, and the extinction rate  $\mu$  was constant throughout time and all lineages.

The simulations were started with a single species. At each time step, each species living in the clade had a probability equal to  $\mu$  to die. If it survived, each species had then a probability given by  $\lambda$  to generate two daughter-species,

otherwise it simply survived to the next time step. This evolution process was simulated during 100 time-steps.

The evolution of  $x$  was so that at each time-step  $x_{t+1} = x_t + \epsilon$  with  $\epsilon \sim N(0, \sigma^2)$ . All values of  $\epsilon$  were independent among all lineages and throughout time. Such a model may not describe properly the evolution of many biological characters, but it was chosen because it can be simulated easily, and it conforms to the assumption of linear change along a branch of the tree.

The simulation parameters were given the following values:  $\alpha = -3$ ,  $\beta = \{0, -1, -2, -3, -4\}$ ,  $\sigma = \{0.01, 0.02\}$ , and  $\mu = \{0, 0.0005, 0.001\}$ . For each combination of parameters the simulation was replicated 100 times. Note that when  $\beta = 0$ , the simulated process was a homogeneous birth–death process with  $\lambda = 1/(1 + e^3) = 0.047$ . Fig. 1 illustrates how the speciation rate varied with respect to  $x$  for the values of  $\beta$  used in these simulations.

Fig. 1

Since the simulations were fully stochastic, so were the number of species living at the end of the simulation. The simulations with 0, 1, or 2 species at the end of the simulation were not considered. In some replications, the number of species were very large (up to 6,807,203 in the present study) which raised some problems because of the data manipulation subsequently needed. I thus considered only the trees with at most 2000 species. When a tree did not meet the selection criterion (between 3 and 2000 species), the simulation was repeated with the same parameter values until 100 trees were obtained for further analyses.

At the end of each simulation, the tree (topology with branch lengths) considering only the species living at the end of the simulation was output together with the values of  $x$  at its tips and nodes. Note that in real applications, the values at the nodes must be reconstructed separately. These data were analyzed with the method described in this paper. The estimates  $\hat{\beta}$ ,

$\hat{\alpha}$ , their standard-errors, the  $P$ -value of the likelihood ratio test of the hypothesis  $\beta = 0$ , the number of species living at the end of the simulation, and the depth of the tree (age of the most recent common ancestor of all living species) were stored for further analyses.

All analyses were done with R version 1.8.1 (R Development Core Team 2003). The simulations were done with a program written in C specially for the present study which was dynamically linked to R. The data analyses were performed with APE, a package written in R for phylogenetic analyses (Paradis et al. 2004).

### *Primates Data*

A complete phylogeny of Primates has been recently reconstructed with the method of supertrees (Purvis 1995). The phylogeny is ultrametric, and has branch lengths in million years (Fig. 2). The tree had several multichotomies (191 nodes for 259 tips): it was necessary to accomodate this since the model used here considers only dichotomous trees. The simplest solution to this problem, and the one adopted here, is to treat these multichotomies as a series of dichotomies with zero-lengthed branches (Purvis and Garland 1993).

Fig. 2

It was interesting to assess whether some biological traits could explain the variation in speciation rate in this order. An obvious candidate was body mass which relationship with diversification rate of Primates has been studied with another approach (Gittleman and Purvis 1998). Furthermore, one of the critical parameters that control speciation is the rate of fixation of new genetic mutations, and this rate is expected to be critically influenced by demography (Gavrilets et al. 2000). Small-bodied species have generally demographic features that make them likely to have a higher rate of fixation than large-bodied ones, such as largely fluctuating numbers (Krebs and Myers 1974),

short dispersal distances (Paradis et al. 1998), or short generation times (Peters 1983). It is therefore logical to expect a negative effect of body size on speciation rate.

Data on body mass were taken from Smith et al. (2003). For the species with no data on body mass in this reference, the mean of the other species belonging to the same genus was calculated: this was done for 43 species. Body mass was still missing for three species which genus is monospecific (*Cebuella pygmaea*, *Simias concolor*, and *Homo sapiens*): data were found for these species in Nowak (1991).

The values of body mass at the nodes of the phylogenies were estimated assuming a Brownian model of trait evolution (Felsenstein 1985). Under this model of trait evolution, the assumption in equation 9 is valid.

## RESULTS

### *Simulation Analysis*

In eight cases out of the 3000 simulations the fitting algorithm did not converge which was clearly due to a small sample size since there were three species in these eight simulated trees.

The overall type I error rate at the nominal level of 5% calculated over all simulations with  $\beta = 0$  was 0.022. For each combination of  $\sigma$  and  $\mu$ , the estimated type I error rates ranged between 0.01 and 0.06. Given that these estimates are based on only 100 replications, there is no evidence that the probability of rejecting the null hypothesis when it is true overrates the nominal level of 5%, or that it is substantially affected by the extinction rate being different from zero.

The estimated power increased with greater values of  $\beta$ , larger values of  $\sigma$ , and smaller values of  $\mu$  (Table 1). When  $\beta$  was large, the effect of the species

Table 1

trait on speciation rate was stronger, and the likelihood ratio test detected this effect with a higher frequency. When  $\sigma$  was high the variation in the species trait was larger, and thus the relation between this trait and speciation rate was tighter and found statistically significant with a higher proportion. When  $\mu > 0$ , the assumptions of the method were not fulfilled, and it may be expected that its statistical properties are affected. This was the case here, but this result depended on the values of the other parameters of the simulations. For instance, the likelihood ratio test was less powerful when  $\beta = -1$  and  $\mu = 0$ , than when  $\beta = -4$  and  $\mu = 0.001$ , even though the assumptions of the method were met in the former case.

Fig. 3

An effect of sample size (number of species) was observed only when  $\mu = 0$  (Fig. 3). As could have been expected, the power of the test was higher for the largest sample sizes. The relation was loose for  $\beta = -1$ , whereas the power was high ( $\approx 1$ ) for  $\beta \leq -2$  when there was 150 or more species. It cannot be excluded that a relationship between power and sample size exists beyond the limit fixed in this study (2000) in the cases with  $\mu > 0$ .

Fig. 4

When considering the precision of the estimators, the results were quite sensitive to the depth of the tree: trees with a most recent common ancestor younger than 60 gave highly dispersed estimates of both parameters. Thus only trees with a depth of at least 60 were considered in the present analysis. The estimates of  $\beta$  were nearly unbiased when  $\mu = 0$ , whereas a bias was observed when  $\mu > 0$ , particularly when  $\sigma = 0.01$  (Fig. 4). This bias was stronger for the larger values of  $\beta$ . On the other hand, the estimates of  $\alpha$  were nearly unbiased in all situations and their distribution showed no variation with respect to  $\beta$  or  $\sigma$ . The median of this distribution was  $-3.04$  (first and third quartiles:  $-3.17$ , and  $-2.92$ ; 98% of the values were distributed between  $-4$  and  $-2$ ).

### *Primates*

Fitting the Yule model (i.e. the null model where  $\lambda$  is constant) gave a deviance of 95.20, and an estimate  $\hat{\lambda} = 0.146$  (SE = 0.009). Fitting a model with an effect of  $\ln(\text{body mass})$  on  $\lambda$  resulted in a deviance of 76.51. The likelihood ratio test of the effect of  $\ln(\text{body mass})$  was statistically significant:  $\chi_1^2 = 18.69$ ,  $P < 0.0001$ . Surprisingly, the effect of  $\ln(\text{body mass})$  was positive:  $\hat{\beta} = 0.21$  (SE = 0.03).

Previous analyses showed a strong contrast between different clades of Primates in terms of diversification rates (Purvis et al. 1995; Paradis 1998a). Particularly, Old World monkeys (Catarrhini) appeared to have diversified at a much higher rate than the other groups of Primates. Given that Old World monkeys have on average larger body masses than New World ones (Platyrrhini), the positive effect of body mass on speciation rate could be due to an artefact. To assess this hypothesis, I fitted a model with a ‘clade’ effect where clade was a categorical variable with five categories: Ape, Catarrhini, Platyrrhini, Strepsirhini (Malagasy lemurs), and *Tarsius*. The states of this variable on the basal nodes of the phylogeny were reconstructed using a parsimony criterion. The deviance of this model was 11.46, and the likelihood ratio test comparing it with the null (Yule) model was  $\chi_4^2 = 83.74$ ,  $P = 0$ . Thus, the improvement in the model fit was dramatic.

I further tested whether  $\ln(\text{body mass})$  had a possible effect on this ‘clade’ effect model. The deviance of this ‘clade +  $\ln(\text{body mass})$ ’ model was 7.21. The likelihood ratio test comparing this model with the ‘clade’ model was slightly significant:  $\chi_1^2 = 4.25$ ,  $P = 0.039$ . Interestingly, the effect of  $\ln(\text{body mass})$  was now negative (Table 2). I finally selected this last model for parameter

estimation. Using equation 4 it is possible to compute the predicted values of speciation rates according to the selected model (Fig. 5). Considering the

Table 2

Fig. 5



observed distributions of body mass in each clade of Primates, the ‘clade’ effect now explains the overall positive relationship between body mass and speciation rate. However, within each clade there is a negative effect of body mass on speciation rate.

A plot of the diagnostics of model fit shows that the terminal branches were particularly well predicted by the selected model (Fig. 2). Long or basal branches were less well predicted, though most branches had a probability greater than or equal to 0.2.

## DISCUSSION

Recent theoretical works on speciation and diversification have been taken in several directions such as characterizing the factors at the origin of speciation (Kaneko and Yomo 2000; Gavrillets 2003; Hochberg et al. 2003; Streelman and Danley 2003), the patterns of phylogenetic diversity (Aldous 1995, 2001; Losos and Adler 1995), or the evolution of genomic complexity (Lynch and Conery 2003). With the increasing quantity of various kinds of data (DNA sequences, ecological and biological data, . . . ), it is crucial to have some statistical methods to test the above theories. The goal of the method presented in this paper is to contribute to fill this need.

There are two main assumptions in the present method: extinction has been absent, and the fully resolved (species-level) phylogeny and the values of the traits under consideration at the tips and the nodes of the phylogeny are all known. These assumptions are discussed in the next four paragraphs.

The model used here (Yule model with covariates) assumes there is no extinction. This assumption is unlikely to be true in many situations as clearly shown by the fossil record of many taxonomic groups. The reason for this assumption is that it makes possible the development of the likelihood function

presented above. It cannot be excluded that a likelihood function may be derived taking extinction into account, but this still needs to be found. Interestingly, the simulation study showed that the type I error rate of the present method is not affected by random extinction: the probability of rejecting the null hypothesis when it was true was not increased when there were extinctions occurring at a constant rate. This suggests that the present method is unlikely to reveal false significant effect of a trait on diversification rate. However, we cannot generalize this result to the cases where extinction rate is heterogeneous. Particularly, it would be interesting to assess whether an extinction rate varying with respect to a species trait (e.g. in the same way modeled here for speciation rate) would be detected by the present method as a significant effect on speciation rate because of the ‘signal’ left on the tree by this heterogeneous extinction rate. This needs to be studied with more extensive simulations than done here.

The power of the test (probability of rejecting the null hypothesis when it is false) was affected by the value of extinction rate: the greater the latter, the smaller the former. However, it should be noted that this depended on the strength of the effect of the trait on the speciation rate: the test of the null hypothesis had a significant power when this effect was strong even though there were extinctions. It should be noted that the simulations were run with a relatively short timespan (100 time steps) which may explain the relative poor performance of the likelihood ratio test in this situation: it was observed that a significant effect was found mostly when  $\beta = -4$ . However, the analysis of the Primates data showed statistically significant estimated coefficients which were much smaller than four in absolute value (as the preliminary analyses of other data did). The fact that the depth of the tree was important for the results of the simulation suggests that the discreteness of the time scale may be critical

here. Ideally, it would be needed to run some simulations with a more realistic time scale ( $> 1000$ ) but these may be computationally hard to tract. Though the estimated power of the tests may not be generalized to real applications, it seems reasonable to assume that the patterns revealed by the simulations (increase in power with sample size, variation in the trait, and strength of the effect) may be generalized.

On a more general note, the estimation of extinction rates with phylogenies of recent species is a subject that should retain more attention. Kubo and Iwasa (1995) showed that the variance of the estimator of extinction rate obtained from molecular phylogenies is too large to be reliable. Paradis (2004) showed that the estimates of extinction rates are negatively biased in a wide range of situations, and that correct estimation of this parameter is possible only when the reconstructed tree is close to the true historical tree of the clade. These results are intuitive since no extinction events are observed in a phylogeny of recent species. Most of the information in this kind of data relates to speciation rate. This justifies, at least partially, to focus the attention on speciation rate. I even suspect that it is not feasible to model extinction rates in the way done here for speciation rates. An additional consideration is that the extinct lineages may have parameters different from those of the extant ones such as a higher extinction rate. Thus, in addition to the problem that these lineages are not observed, it is possible that there is a problem of heterogeneity in the parameters as well.

The present method requires to know the fully resolved phylogeny of the species, and the values of the traits for these species at the tips and at the nodes of the phylogeny. It is generally not possible to have direct observations of the traits for the nodes of the trees, so they must be estimated using ancestral character reconstruction methods (Webster and Purvis 2002, and

references therein). The assumption of a fully resolved tree can also be relaxed since a tree with multichotomies can be analyzed. Ideally, one needs to take such uncertainty into account in the estimation procedure so that the standard-errors of the parameters may be corrected. One possible approach is to repeat the analysis with different trees and different values of the ancestral traits, and then quantify the resulting variability on the parameter estimates as a component of their standard-errors due to this uncertainty (Anderson et al. 2001). Another, more sophisticated, approach is to use directly the confidence intervals on trees and ancestral trait values (see Schluter et al. 1997) to compute a component of the variance of the estimators due to model selection uncertainty (Buckland et al. 1997; Burnham and Anderson 2002).

The method in the form presented here considers only linear effects. Though this gives many possibilities such as including several variables and possible interactions between them, it is possible that nonlinear models may be more biologically realistic models of the effects of species traits on speciation rates. Such candidate models are threshold models where the effect of a variable differs depending on some threshold values. It is straightforward to extend the present method to nonlinear models: the only condition is that it must be possible to derive the value of speciation rate with respect to the variables included in the model (such as equation 4).

Throughout this paper, it has been assumed that continuous traits evolved in a simple way (Brownian motion), but another model of change may be used provided that the integration along each branch may be expressed in a way that it can be incorporated in the likelihood function. It should be noted that it is not required to know the exact temporal variation in species traits and speciation rate, but only the integral of the latter along each branch of the tree. For instance, it does not matter whether the speciation rate changed linearly or

with some random fluctuations along a branch, the integral will be the same in both cases.

The present method gives some emphasis on parameter estimation by contrast to hypothesis testing. With this respect, it is in agreement with some recent trends in statistical science and data analysis (Johnson 1995, 1999; Nelder 1999; Anderson et al. 2000; Venables and Ripley 2002). In evolutionary biology, null hypothesis testing is still extremely popular. I believe much could be gained from giving more emphasis on parameter estimation in this field. This would allow the development of predictive models *from* the data (i.e. with an empirical ground). Null hypotheses usually focus on simple models whereas more complex models are often needed to describe biological processes. Another advantage of parameter estimation approaches is that the estimates (under some obvious conditions) are comparable among studies. Considering more carefully parameter estimation does not deny the importance of statistical tests as illustrated here where likelihood ratio tests were used to select models with the appropriate variables.

The results presented in this paper are an illustration of the potentialities of the present method, but they clearly need some comments. The data analyzed here were nearly the same than in Gittleman and Purvis (1998). With respect to body size, the results from both studies were in the same direction: a positive effect on diversification was found for primates; however, the studies disagree with respect to the significance of this effect. This may come either from the different methods used, or from differences in the data. Gittleman and Purvis (1998) used the phylogenetically-based contrasts method (Barraclough et al. 1998) which is likely to explain the discrepancy with the present study.

An important issue raised by the Primates example is the treatment of multichotomies by the present method. The Yule model (with or without

covariates) assumes that all speciation events are dichotomous, and thus that a phylogenetic tree under analysis should be fully dichotomous as well. However, many phylogenies at the species level have multichotomies. If the latter are the results of a series of speciation events so close in time that they cannot be resolved in the reconstructed phylogeny, then it is valid to treat them as a series of dichotomies with zero-lengthed branches. However, if the multichotomy is a consequence of a lack of knowledge (a so-called “soft polytomy”), then resolving the multichotomy in this way may be more problematic. The crucial point is that ancestral traits should be correctly estimated. Several studies suggest that unresolved multichotomies do not bias the estimation of ancestral traits (Felsenstein 1985; Purvis and Garland 1993; Garland and Díaz-Uriarte 1999). On the other hand, errors in the topology of the tree is a more serious problem for the estimation of ancestral trait values (Symonds 2002). It seems that unresolved multichotomies will tend to hide actual relationships, and thus increase type II error rates, rather than increase type I error rates (Garland and Díaz-Uriarte 1999; Symonds 2002). This obviously needs further study.

It has been possible to show that the apparent overall positive effect of body mass on speciation rate was actually an artefact due to differences among the different clades of Primates. The most rapidly diversifying clades (Cathartini, Ape) have, on average, larger body sizes than the slowly diversifying clades (*Tarsius*, Strepsirhini): this created an apparent positive relationship between body size and speciation rate when the clade-effect was not taken into account. This among-clade heterogeneity in diversification has already been characterized in previous studies (Purvis et al. 1995; Paradis 1998a). Similar clade-specific diversification has been shown in Hawaiian birds (Lovette et al. 2002). An interesting result from the present analysis is that, after taking into account this among-clade heterogeneity, the effect of body mass on speciation rate was

negative, though only slightly statistically significant. It is possible that the effect of body mass, if any, on speciation rate is too weak in this mammal order to be characterized clearly, as was suggested by Gittleman and Purvis (1998). Furthermore, recent assessments of the effect of body size on species diversity on a large scale have not given clear-cut results (Orme et al. 2002a,b).

To conclude, I believe the present method will bring the opportunity to assess many hypotheses about macroevolutionary processes. It will hopefully fill a gap between the approach of estimating speciation and extinction rates and biological and ecological hypotheses on diversification.

#### ACKNOWLEDGMENTS

I am grateful to two anonymous referees and the Associate Editor for their constructive comments on a previous version of this paper. This research was financially supported by the Institut Français de la Biodiversité and the Centre National de la Recherche Scientifique. This is publication 2004-XXX of the Institut des Sciences de l'Évolution (Unité Mixte de Recherche 5554 du Centre National de la Recherche Scientifique).

#### LITERATURE CITED

- Aldous, D. 1995. Darwin's log: a toy model of speciation and extinction. *J. Appl. Prob.* 32:279–295.
- Aldous, D. J. 2001. Stochastic models and descriptive statistics for phylogenetic trees, from Yule to today. *Statist. Sci.* 16:23–34.
- Anderson, D. R., K. P. Burnham, and W. L. Thompson. 2000. Null hypothesis testing: problems, prevalence, and an alternative. *J. Wildl. Manage.* 64:912–923.
- Anderson, D. R., W. A. Link, D. H. Johnson, and K. P. Burnham. 2001.

- Suggestions for presenting the results of data analyses. *J. Wildl. Manage.* 65:373–378.
- Archibald, J. D. and D. H. Deutschman. 2001. Quantitative analysis of the timing of the origin and diversification of extant placental orders. *J. Mamm. Evol.* 8:107–124.
- Barracough, T. G. and S. Nee. 2001. Phylogenetics and speciation. *Trends Ecol. Evol.* 16:391–399.
- Barracough, T. G., S. Nee, and P. H. Harvey. 1998. Sister-group analysis in identifying correlates of diversification. *Evol. Ecol.* 12:751–754.
- Benton, M. J., M. A. Wills, and R. Hitchin. 2000. Quality of the fossil record through time. *Nature* 403:534–537.
- Bleiweiss, R. 1998. Fossil gap analysis supports early Tertiary origin of trophically diverse avian orders. *Geology* 26:323–326.
- Bokma, F. 2002. A statistical test of unbiased evolution of body size in birds. *Ecology* 56:2499–2504.
- Bokma, F. 2003. Testing for equal rates of cladogenesis in diverse taxa. *Ecology* 57:2469–2474.
- Buckland, S. T., K. P. Burnham, and N. H. Augustin. 1997. Model selection: an integral part of inference. *Biometrics* 53:603–618.
- Burnham, K. P. and D. R. Anderson. 2002. Model selection and multimodel inference. A practical information-theoretic approach (second edition). Springer, New York.
- Cardillo, M., J. S. Huxtable, and L. Bromham. 2003. Geographic range size, life history and rates of diversification in Australian mammals. *J. Evol. Biol.* 16:282–288.
- Cox, D. R. and D. Oakes. 1984. Analysis of survival data. Monographs on statistics and applied probability, Chapman and Hall, London.



- Darwin, J. H. 1956. The behaviour of an estimator for a simple birth and death process. *Biometrika* 43:23–31.
- Felsenstein, J. 1985. Phylogenies and the comparative method. *Am. Nat.* 125:1–15.
- Fisher, D. O., S. P. Blomberg, and I. P. F. Owens. 2003. Extrinsic versus intrinsic factors in the decline and extinction of Australian marsupials. *Proc. R. Soc. Lond. B* 270:1801–1808.
- Foote, M., J. P. Hunter, C. M. Janis, and J. J. Sepkoski, Jr. 1999. Evolutionary and preservational constraints on origins of biologic groups: divergence times of eutherian mammals. *Science* 283:1310–1314.
- Foote, M. and J. J. Sepkoski, Jr. 1999. Absolute measures of the completeness of the fossil record. *Nature* 398:415–417.
- Freedman, D. 1995/96. Some issues in the foundations of statistics (with discussion). *Foundations of Science* 1:19–83.
- Garland, T., Jr. and R. Díaz-Uriarte. 1999. Polytomies and phylogenetically independent contrasts: examination of the bounded degrees of freedom approach. *Syst. Biol.* 48:547–558.
- Gavrilets, S. 2003. Models of speciation: what have we learned in 40 years? *Evolution* 57:2197–2215.
- Gavrilets, S., R. Acton, and J. Gravner. 2000. Dynamics of speciation and diversification in a metapopulation. *Evolution* 54:1493–1501.
- Gittleman, J. L. and A. Purvis. 1998. Body size and species-richness in carnivores and primates. *Proc. R. Soc. Lond. B* 265:113–119.
- Hochberg, M. E., B. Sinervo, and S. P. Brown. 2003. Socially mediated speciation. *Evolution* 57:154–158.
- Hunter, J. P. 1998. Key innovations and the ecology of macroevolution. *Trends Ecol. Evol.* 13:31–36.

- Ihaka, R. and R. Gentleman. 1996. R: a language for data analysis and graphics. *J. Comput. Graph. Statist.* 5:299–314.
- Isaac, N. J. B., P. M. Agapow, P. H. Harvey, and A. Purvis. 2003. Phylogenetically nested comparisons for testing correlates of species richness: a simulation study of continuous variables. *Evolution* 57:18–26.
- Jablonski, D. 2000. Micro- and macroevolution: scale and hierarchy in evolutionary biology and paleobiology. *Paleobiology* 26:15–52.
- Johnson, D. H. 1995. Statistical sirens: the allure of nonparametrics. *Ecology* 76:1998–2000.
- Johnson, D. H. 1999. The insignificance of statistical significance testing. *J. Wildl. Manage.* 63:763–772.
- Kaneko, K. and T. Yomo. 2000. Sympatric speciation: compliance with phenotype diversification from a single genotype. *Proc. R. Soc. Lond. B* 267:2367–2373.
- Krebs, C. J. and J. H. Myers. 1974. Population cycles in small mammals. *Advances in Ecological Research* 8:267–399.
- Kubo, T. and Y. Iwasa. 1995. Inferring the rates of branching and extinction from molecular phylogenies. *Evolution* 49:694–704.
- Losos, J. B. and F. R. Adler. 1995. Stumped by trees? A generalized null model for patterns of organismal diversity. *Am. Nat.* 145:329–342.
- Lovette, I. J., E. Bermingham, and R. E. Ricklefs. 2002. Clade-specific morphological diversification and adaptive radiation in Hawaiian songbirds. *Proc. R. Soc. Lond. B* 269:37–42.
- Lynch, M. and J. S. Conery. 2003. The origins of genome complexity. *Science* 302:1401–1404.
- Marshall, C. R. 1997. Confidence intervals on stratigraphic ranges with nonrandom distributions of fossil horizons. *Paleobiology* 23:165–173.

- McCullagh, P. and J. A. Nelder. 1989. Generalized linear models (second edition). Chapman & Hall, London.
- Mooers, A. Ø. and S. B. Heard. 1997. Evolutionary process from phylogenetic tree shape. *Quart. Rev. Biol.* 72:31–54.
- Nee, S. 2001. Inferring speciation rates from phylogenies. *Evolution* 55:661–668.
- Nee, S., R. M. May, and P. H. Harvey. 1994. The reconstructed evolutionary process. *Phil. Trans. R. Soc. Lond. B* 344:305–311.
- Nelder, J. A. 1999. From statistics to statistical science (with discussion). *Statistician* 48:257–269.
- Nichols, J. D., R. W. Morris, C. Brownie, and K. H. Pollock. 1986. Sources of variation in extinction rates, turnover, and diversity of marine invertebrate families during the Paleozoic. *Paleobiology* 12:421–432.
- Nowak, R. M. 1991. Walker's mammals of the world (fifth edition). The John Hopkins University Press, Baltimore and London.
- Orme, C. D. L., N. J. B. Isaac, and A. Purvis. 2002a. Are most species small? Not within species-level phylogenies. *Proc. R. Soc. Lond. B* 269:1279–1287.
- Orme, C. D. L., D. L. J. Quicke, J. M. Cook, and A. Purvis. 2002b. Body size does not predict species richness among the metazoan phyla. *J. Evol. Biol.* 15:235–247.
- Paradis, E. 1997. Assessing temporal variations in diversification rates from phylogenies: estimation and hypothesis testing. *Proc. R. Soc. Lond. B* 264:1141–1147.
- Paradis, E. 1998a. Detecting shifts in diversification rates without fossils. *Am. Nat.* 152:176–187.
- Paradis, E. 1998b. Testing for constant diversification rates using molecular phylogenies: a general approach based on statistical tests for goodness of fit. *Mol. Biol. Evol.* 15:476–479.

- Paradis, E. 2003. Analysis of diversification: combining phylogenetic and taxonomic data. *Proc. R. Soc. Lond. B* 270:2499–2505.
- Paradis, E. 2004. Can extinction rates be estimated without fossils? *J. Theor. Biol.* 229:19–30.
- Paradis, E., S. R. Baillie, W. J. Sutherland, and R. D. Gregory. 1998. Patterns of natal and breeding dispersal in birds. *J. Anim. Ecol.* 67:518–536.
- Paradis, E., J. Claude, and K. Strimmer. 2004. APE: analyses of phylogenetics and evolution in R language. *Bioinformatics* 20:289–290.
- Peters, R. H. 1983. The ecological implications of body size. Cambridge Series in Ecology, Cambridge University Press, Cambridge.
- Purvis, A. 1995. A composite estimate of primate phylogeny. *Phil. Trans. R. Soc. Lond. B* 348:405–421.
- Purvis, A. and T. Garland, Jr. 1993. Polytomies in comparative analyses of continuous characters. *Syst. Biol.* 42:569–575.
- Purvis, A., S. Nee, and P. H. Harvey. 1995. Macroevolutionary inferences from primate phylogeny. *Proc. R. Soc. Lond. B* 260:329–333.
- Pybus, O. G. and P. H. Harvey. 2000. Testing macro-evolutionary models using incomplete molecular phylogenies. *Proc. R. Soc. Lond. B* 267:2267–2272.
- R Development Core Team, 2003. R: a language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL: <http://www.R-project.org>.
- Sanderson, M. J. and G. Bharathan. 1993. Does cladistic information affect inferences about branching rates? *Syst. Biol.* 42:1–17.
- Sanderson, M. J. and M. J. Donoghue. 1994. Shifts in diversification rate with the origin of angiosperms. *Science* 264:1590–1593.
- Schluter, D., T. Price, A. Ø. Mooers, and D. Ludwig. 1997. Likelihood of ancestor states in adaptive radiation. *Evolution* 51:1699–1711.

- Schnabel, R. B., J. E. Koontz, and B. E. Weiss. 1985. A modular system of algorithms for unconstrained minimization. *ACM Trans. Math. Software* 11:419–440.
- Slowinski, J. B. and C. Guyer. 1993. Testing whether certain traits have caused amplified diversification: an improved method based on a model of random speciation and extinction. *Am. Nat.* 142:1019–1024.
- Smith, F. A., S. K. Lyons, S. K. M. Ernest, K. E. Jones, D. M. Kaufman, T. Dayan, P. A. Marquet, J. H. Brown, and J. P. Haskell. 2003. Body mass of late quaternary mammals. *Ecology* 84:3403. URL: <http://www.esapubs.org/archive/ecol/E084/094/default.htm>.
- Streelman, J. T. and P. D. Danley. 2003. The stages of vertebrate evolutionary radiation. *Trends Ecol. Evol.* 18:126–131.
- Symonds, M. R. E. 2002. The effects of topological inaccuracy in evolutionary trees on the phylogenetic comparative method of independent contrasts. *Syst. Biol.* 51:541–553.
- Tipper, J. C. 1998. The influence of field sampling area on estimates of stratigraphic completeness. *J. Geol.* 106:727–739.
- Venables, W. N. and B. D. Ripley. 2002. *Modern applied statistics with S* (fourth edition). Springer, New York.
- Webster, A. J. and A. Purvis. 2002. Testing the accuracy of methods for reconstructing ancestral states of continuous characters. *Proc. R. Soc. Lond. B* 269:143–149.
- Weiss, R. E. and C. R. Marshall. 1999. The uncertainty in the true end point of a fossil's stratigraphic range when stratigraphic sections are sampled discretely. *Math. Geol.* 31:435–453.
- Wollenberg, K., J. Arnold, and J. C. Avise. 1996. Recognizing the forest for the trees: testing temporal patterns of cladogenesis using a null model of

stochastic diversification. Mol. Biol. Evol. 13:833–849.

## APPENDIX

Some details are given here on how discrete (categorical) variables are handled in the present method, and how the corresponding parameters must be interpreted. For instance, consider color as such a variable which could take two values: red or green. To include this variable in a model, it will be substituted by a numeric in the following way:

$$\begin{aligned}\text{red} &\rightarrow 0 \\ \text{green} &\rightarrow 1\end{aligned}$$

The model fit will give two parameter estimates:  $\hat{\beta}$  and  $\hat{\alpha}$ . The predicted values of speciation rate are then obtained with:

$$\begin{aligned}\text{red:} \quad \text{logit}(\lambda) &= \hat{\alpha} \\ \text{green:} \quad \text{logit}(\lambda) &= \hat{\beta} + \hat{\alpha}\end{aligned}$$

where  $\text{logit}(\lambda) = \ln[\lambda/(1 - \lambda)]$ .

In the general case of a categorical variable with  $n$  categories,  $n - 1$  binary variables are created. For instance, with the primate data above, the four numeric variables were:

$$\begin{aligned}\text{Ape} &\rightarrow 0 \ 0 \ 0 \ 0 \\ \text{Catarrhini} &\rightarrow 1 \ 0 \ 0 \ 0 \\ \text{Platyrrhini} &\rightarrow 0 \ 1 \ 0 \ 0 \\ \text{Strepsirhini} &\rightarrow 0 \ 0 \ 1 \ 0 \\ \text{\textit{Tarsius}} &\rightarrow 0 \ 0 \ 0 \ 1\end{aligned}$$

This led to the estimation of five parameters:  $\hat{\beta}_1$ ,  $\hat{\beta}_2$ ,  $\hat{\beta}_3$ ,  $\hat{\beta}_4$ , and  $\hat{\alpha}$ . The predicted values of speciation rate were obtained with:

$$\begin{aligned}
\text{Ape:} & \quad \text{logit}(\lambda) = \hat{\alpha} \\
\text{Catarrhini:} & \quad \text{logit}(\lambda) = \hat{\beta}_1 + \hat{\alpha} \\
\text{Platyrrhini:} & \quad \text{logit}(\lambda) = \hat{\beta}_2 + \hat{\alpha} \\
\text{Strepsirhini:} & \quad \text{logit}(\lambda) = \hat{\beta}_3 + \hat{\alpha} \\
\text{\textit{Tarsius}:} & \quad \text{logit}(\lambda) = \hat{\beta}_4 + \hat{\alpha}
\end{aligned}$$

In a model with a categorical variable (as color above) and a continuous variable (say  $x$ ), it is possible to include, in addition to the main effects of both variables, an interaction term which is numerically coded as the product of the numeric variable coding for color and  $x$ :

$$\begin{aligned}
\text{red} & \quad \rightarrow \quad 0 \\
\text{green} & \quad \rightarrow \quad x
\end{aligned}$$

The model fit will give four parameter estimates:  $\hat{\beta}_1$  (for  $x$ ),  $\hat{\beta}_2$  (for color),  $\hat{\beta}_3$  (for the interaction), and  $\hat{\alpha}$ . The predicted values of speciation rate are then obtained with:

$$\begin{aligned}
\text{red:} & \quad \text{logit}(\lambda) = \hat{\beta}_1 x + \hat{\alpha} \\
\text{green:} & \quad \text{logit}(\lambda) = (\hat{\beta}_1 + \hat{\beta}_3)x + \hat{\beta}_2 + \hat{\alpha}
\end{aligned}$$

This shows clearly that the interaction term is interpreted as a contrast in the effect of  $x$  with respect to the different categories of color (if  $\hat{\beta}_3 = 0$ , then only the intercept will be different). In the general case of a categorical variable with  $n$  categories,  $n - 1$  variables will be created with the product of the  $n - 1$  binary variables with  $x$ .

Similarly, a model with two categorical variables may have an interaction term numerically coded by all possible products of the individual numeric variables: a full model including the main and interaction effects of two categorical variables with  $n_1$  and  $n_2$  categories, respectively, will thus have  $n_1 n_2 - 1$  parameters ( $n_1 - 1$  for the main effect of the first variable,  $n_2 - 1$  for the second, and  $(n_1 - 1)(n_2 - 1)$  for the interaction).

FIG. 1. Examples of variation of speciation rate ( $\lambda$ ) with respect to a trait ( $x$ ). These values were used in the simulation study.

FIG. 2. Phylogeny of Primates from Purvis (1995), and barplot of  $\ln(\text{body mass})$  for each species. The grey level indicates the probability that the given branch is observed according to the selected model.

FIG. 3. Relationship between the  $P$ -values of the likelihood ratio test of the hypothesis  $\beta = 0$  and the number of species in the simulated trees for the different values of  $\beta$  and extinction rate ( $\mu$ ).

FIG. 4. Relationship between the true values of  $\beta$  and its estimates  $\hat{\beta}$  for the different values of rate of trait evolution ( $\sigma$ ) and extinction rate ( $\mu$ ). The dashed lines are  $x = y$ . Only the simulated trees with a depth of at least 60 were considered.

FIG. 5. Predicted values of speciation rate ( $\lambda$ ) with respect to body mass for the different clades of primates according to the selected model. The bold parts of the curves show the range of observed values of body mass for each clade. The observed values of body mass for each species are indicated on the inner side of the  $x$ -axis.



Table 1: Estimated statistical power for each combination of strength of effect of species trait on speciation rate ( $\beta$ ), trait evolution rate ( $\sigma$ ), and extinction rate ( $\mu$ ).

| $\beta$ | $\sigma$ | $\mu$ |        |       |
|---------|----------|-------|--------|-------|
|         |          | 0     | 0.0005 | 0.001 |
| -1      | 0.01     | 0.05  | 0.01   | 0.03  |
|         | 0.02     | 0.15  | 0.05   | 0.01  |
| -2      | 0.01     | 0.17  | 0.03   | 0.02  |
|         | 0.02     | 0.49  | 0.10   | 0.06  |
| -3      | 0.01     | 0.38  | 0.10   | 0.04  |
|         | 0.02     | 0.68  | 0.26   | 0.12  |
| -4      | 0.01     | 0.58  | 0.05   | 0.09  |
|         | 0.02     | 0.58  | 0.26   | 0.22  |

Table 2: Parameter estimates for Primates.

| Effect        |                | Estimate | SE   |
|---------------|----------------|----------|------|
| clade         | Ape            | -        | -    |
|               | Catarrhini     | 1.22     | 0.24 |
|               | Platyrrhini    | -0.36    | 0.27 |
|               | Strepsirhini   | -0.86    | 0.30 |
|               | <i>Tarsius</i> | -1.91    | 0.46 |
| ln(body mass) |                | -0.17    | 0.06 |
| intercept     |                | -0.32    | 0.60 |

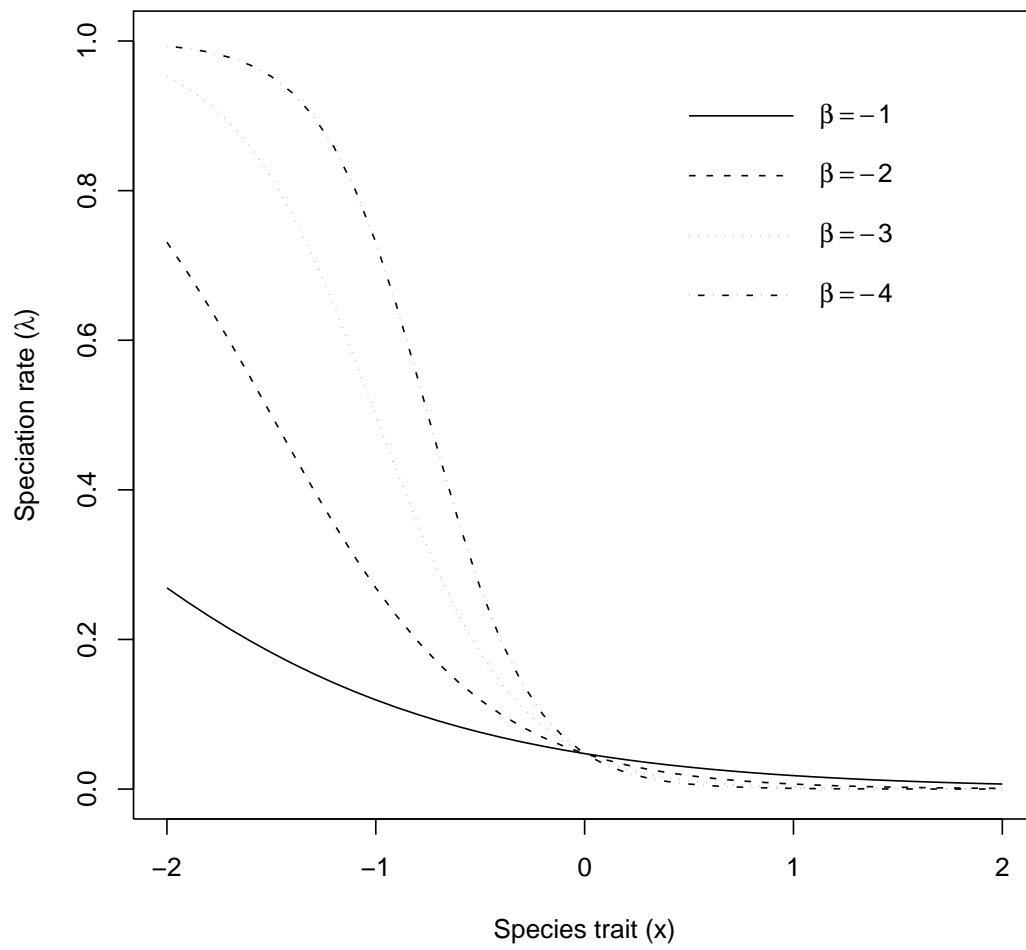


Figure 1:

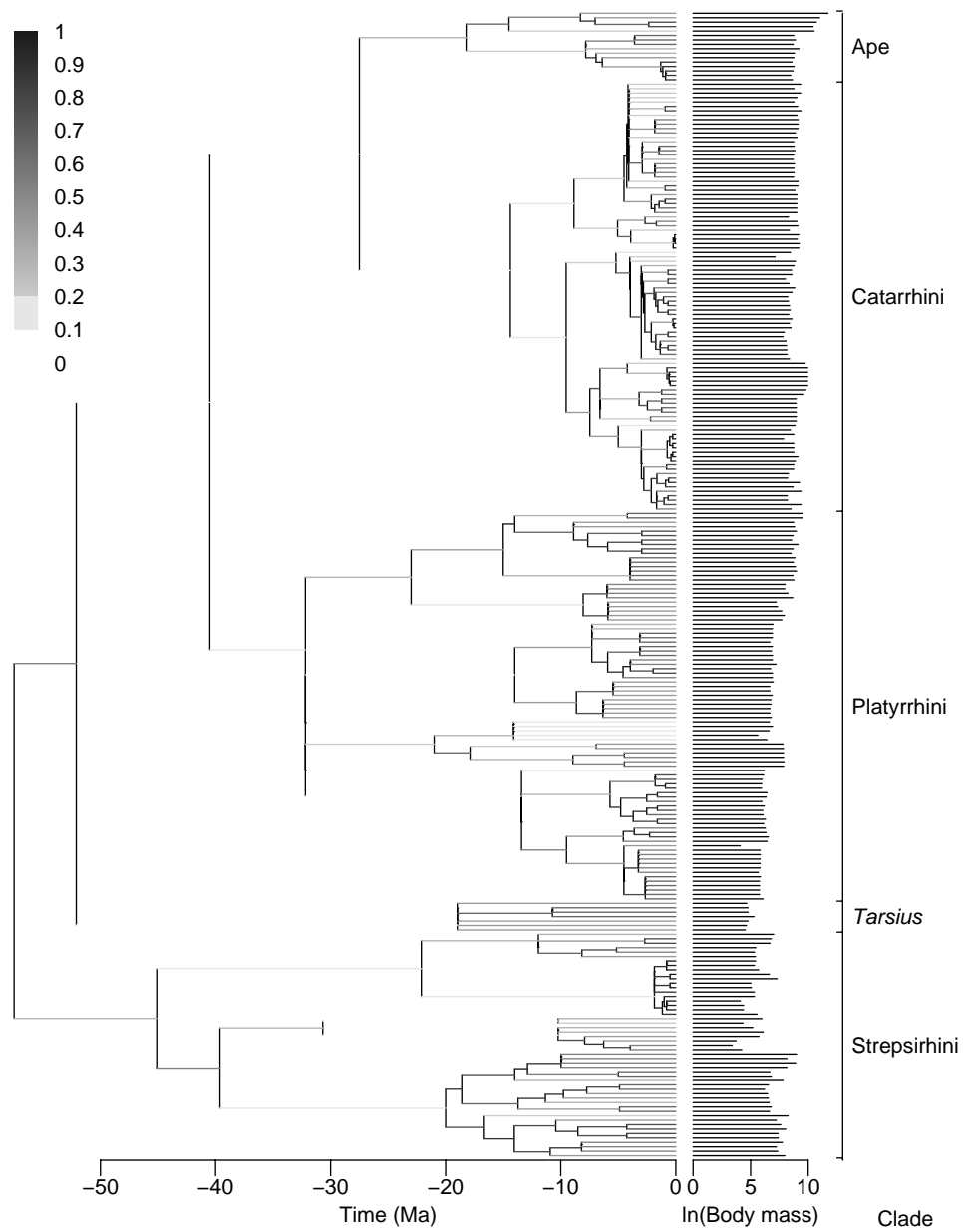


Figure 2:

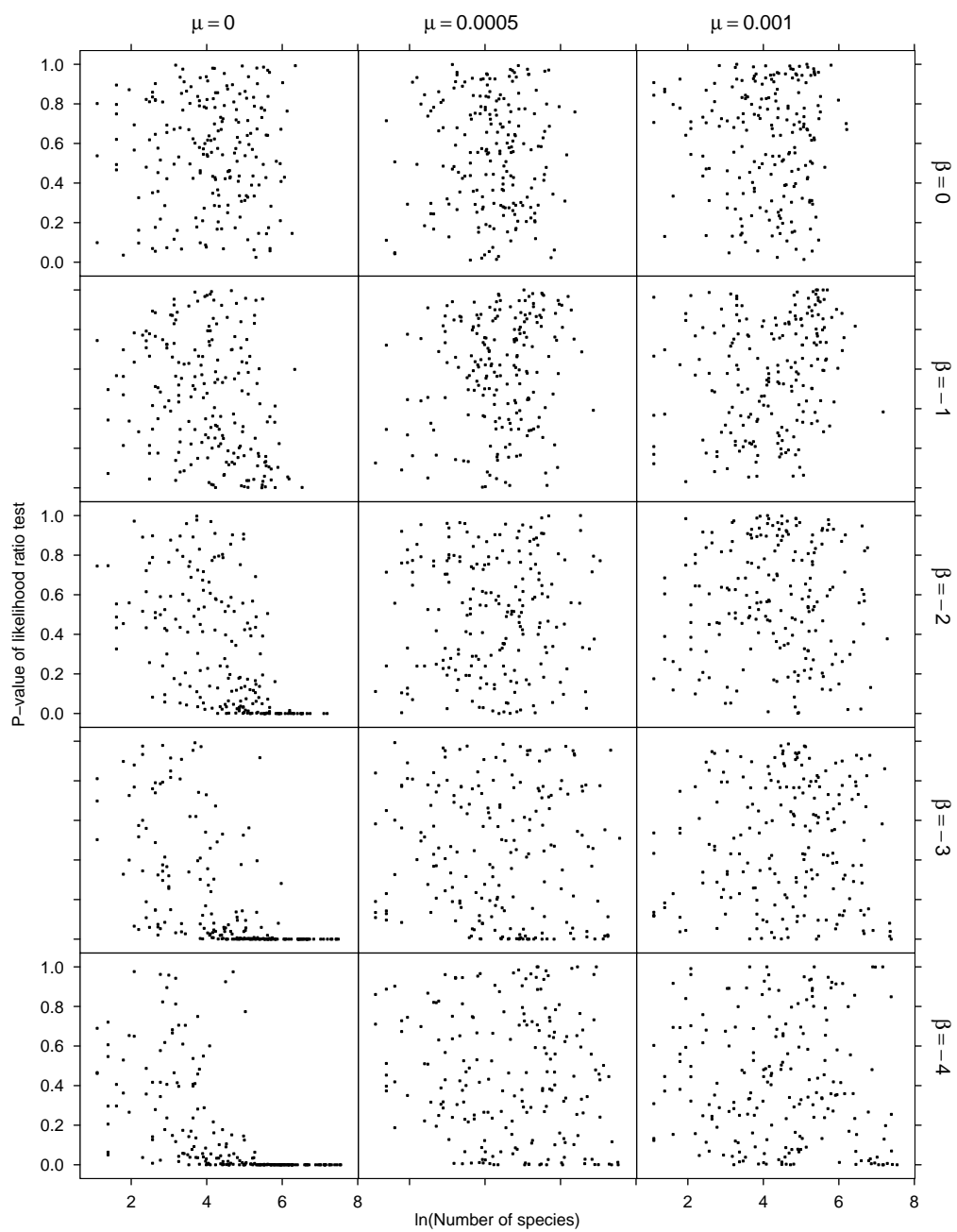


Figure 3:

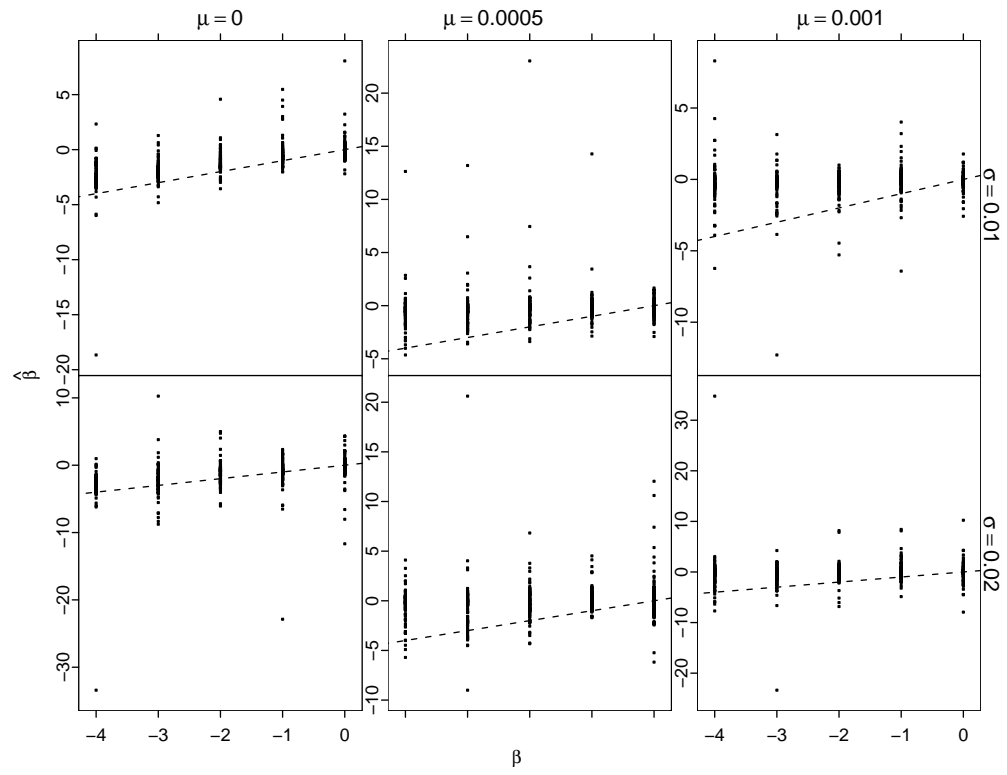


Figure 4:

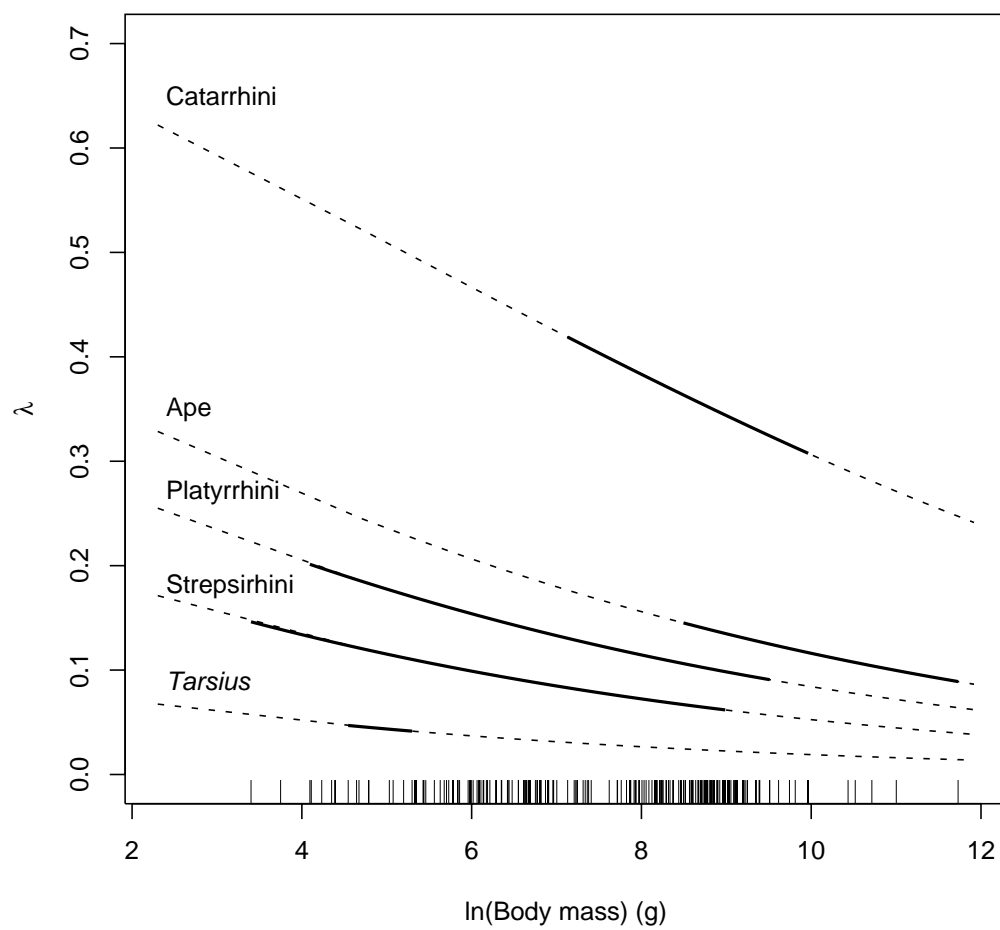


Figure 5: