

WILEY



Confidence Limits on Phylogenies: An Approach Using the Bootstrap

Author(s): Joseph Felsenstein

Source: *Evolution*, Vol. 39, No. 4 (Jul., 1985), pp. 783-791

Published by: Society for the Study of Evolution

Stable URL: <http://www.jstor.org/stable/2408678>

Accessed: 17-11-2016 19:19 UTC

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at
<http://about.jstor.org/terms>



Wiley, Society for the Study of Evolution are collaborating with JSTOR to digitize, preserve and extend access to *Evolution*

CONFIDENCE LIMITS ON PHYLOGENIES: AN APPROACH USING THE BOOTSTRAP

JOSEPH FELSENSTEIN

Department of Genetics SK-50, University of Washington, Seattle, WA 98195

Abstract.—The recently-developed statistical method known as the “bootstrap” can be used to place confidence intervals on phylogenies. It involves resampling points from one’s own data, with replacement, to create a series of bootstrap samples of the same size as the original data. Each of these is analyzed, and the variation among the resulting estimates taken to indicate the size of the error involved in making estimates from the original data. In the case of phylogenies, it is argued that the proper method of resampling is to keep all of the original species while sampling characters with replacement, under the assumption that the characters have been independently drawn by the systematist and have evolved independently. Majority-rule consensus trees can be used to construct a phylogeny showing all of the inferred monophyletic groups that occurred in a majority of the bootstrap samples. If a group shows up 95% of the time or more, the evidence for it is taken to be statistically significant. Existing computer programs can be used to analyze different bootstrap samples by using weights on the characters, the weight of a character being how many times it was drawn in bootstrap sampling. When all characters are perfectly compatible, as envisioned by Hennig, bootstrap sampling becomes unnecessary; the bootstrap method would show significant evidence for a group if it is defined by three or more characters.

Received July 12, 1984. Accepted April 12, 1985

It is rare that any attempt is made to put a confidence interval on an estimate of a phylogeny. Most methods for inferring phylogenies yield one or a few trees, and their users rarely go beyond examining the variation among trees that are tied with the best tree under whatever criterion is being employed. There is no reason to believe that this practice constitutes an adequate exploration of the size of the confidence limits on the estimate.

A few authors have explored the question of confidence limits on phylogenies. The pioneer in doing so is Cavender (1978, 1981) who examined the confidence limits for a four-species case, in terms of how many steps worse a tree must be than the most parsimonious tree to be significantly worse. His results were a bit disconcerting: when inferences were based on twenty characters, a tree would have to be 9 steps worse to be significantly worse. This implies that the confidence intervals would be quite large.

Templeton (1983) has constructed a test of whether one tree is significantly better supported than another. In principle such a test could be used to delimit

a confidence interval by finding all trees that cannot be rejected in comparison with the best supported tree. I have recently extended Cavender’s analysis to the case of a molecular clock with three species, obtaining, in that case, confidence limits that were somewhat smaller than Cavender’s (Felsenstein, 1985). I have also recently reviewed the application of statistics to inferring phylogenies (Felsenstein, 1983a); that paper may be consulted for earlier references on statistical estimation of phylogenies.

An important recent statistical method is the bootstrap (Efron, 1979), a relative of the jackknife. Like the jackknife, it is a method of resampling one’s own data to infer the variability of the estimate. This paper will explore the use of the bootstrap in inferring phylogenies, where it leads to a practical method for placing confidence intervals on the estimates.

The Bootstrap

A straightforward statistical exposition of the bootstrap is given by Efron and Gong (1983), and a readable elementary treatment is that by Diaconis and Efron (1983). The basic idea of the bootstrap

involves inferring the variability in an unknown distribution from which your data were drawn by resampling from the data. Suppose that you had data points x_1, x_2, \dots, x_n , which you are willing to assume were drawn independently from the same distribution. From these, applying some method T of statistical estimation, we obtain an estimate

$$t = T(x_1, x_2, \dots, x_n) \quad (1)$$

of a parameter we are interested in. If we knew the exact distribution from which the x_i were drawn, and if the function T were sufficiently tractable algebraically, we could obtain a formula for the standard error of the estimate t , and also construct confidence intervals for t .

The bootstrap procedure is most useful when we either do not know the distribution of the x_i , or when T is so complicated that its standard error is difficult to compute. It suggests that we resample our data to construct a series of fictional sets of data. Each of these is constructed by sampling n points from the x_i , sampling with replacement. Each such fictional data set consists of n points, x_1^*, \dots, x_n^* where each point x_i^* is drawn at random from among the n original data points. It is quite likely that, in this resampling process, some of the original data points are represented more than once, and others are omitted.

For each fictional set of data, we compute the estimate

$$t^* = T(x_1^*, x_2^*, \dots, x_n^*). \quad (2)$$

The resampling process is done many times (say r times), each time producing a fictional sample of n points by sampling with replacement from the original n data points. For each the estimate t^* is computed. We are then in possession of a collection of r estimates of the parameter. The essential idea of the bootstrap is that this set of estimates has a distribution that approximates the distribution of the actual estimate t . A bias-corrected estimate of the parameter can be computed by averaging the r different t^* values (Efron and Gong, 1983). The variance of t

can be inferred by computing the variance of this collection of t^* values, and the confidence limits on the parameter can be approximated by using the appropriate upper and/or lower percentiles of the observed distribution of the t^* values.

The justification for this resampling is that, if the original sample size n is large, each possible value of x will be represented in the same proportion as in the underlying distribution, and resampling from the data points with replacement will be the same as sampling from the underlying distribution. For smaller sample sizes, the process is an approximation but frequently is a very good one. The monograph by Efron (1982) can be consulted for further details on the properties of the bootstrap.

Bootstrapping Phylogenies

How can the bootstrap be applied to phylogenies? Instead of sample points x_1, x_2, \dots, x_n we usually have a table of species \times characters (or species \times sites for molecular sequences). It is not immediately obvious how resampling can be done in the data table. I will argue that a justifiable procedure is to bootstrap across the characters, that is, to sample characters (or sites) from the data table with replacement. Thus, each bootstrap sample consists of a new data table with the same set of species, but with some of the original characters duplicated and others dropped by the process of sampling n characters from the original set with replacement.

The justification for this is that we can view each character as having evolved independently from the others according to a stochastic process that has among its parameters the topology and branch lengths of the underlying phylogeny. Each character is then a random sample from a distribution of all possible configurations of characters. For example, if we are considering nucleic acid sequence data with p species, there are 4^p possible outcomes at each site, not counting the possibilities of deletion and insertion. To a first approximation we can consider each

site to be independently drawn from a distribution with 4^p possibilities, whose probabilities depend on the phylogeny we are trying to estimate.

Given this independence of evolutionary processes in different characters, the configurations in the characters are seen to be drawn independently and identically distributed (i.i.d.), a necessary condition for the bootstrap method to be valid. In fact, in the case of discrete character states (such as nucleic acids), the underlying distribution is multinomial, since there are 4^p possibilities each of which has some probability of occurring. Despite the complexity of the structure being inferred (the phylogeny) the statistical model is a very straightforward one—independent samples from a multinomial distribution.

It might be argued that this presupposes that the same probabilistic evolutionary process is operating in all of the characters, which is extremely unrealistic. Such an assumption is not necessary. If instead we had a variety of different kinds of characters evolving according to different processes, we need only imagine that there is an additional stage in the process of random sampling, one occurring in the mind of the systematist. We imagine, as part of the stochastic process, a step in which the systematist randomly draws each character from a pool of different kinds of characters, each kind having a different evolutionary process that applies to it. Once drawn, each character then has its actual configuration determined by the appropriate stochastic evolutionary process. The resulting distribution of character configurations is a mixture of multinomial distributions, and, as such, is still a multinomial distribution and is still i.i.d.

In practice the systematist may not have sampled the characters at random. Systematists frequently include characters in the study in groups (such as groups of measurements on the skull). We are then not justified in regarding the process of choice of characters as a series of random samples from a pool of possible

characters. I have recently discussed (Felsenstein, 1983a) some of the statistical issues involved in such a random-sampling model of inference of phylogenies.

A more serious difficulty is lack of independence of the evolutionary processes in different characters. If the characters are correlated (as measurement characters often are), then, in effect, we have fewer characters in the study than we believe. If correlations mean that what appear to be 50 independent characters are really more like 30, the variability that we infer for our estimate will be too small, producing overconfidence in the result; a bootstrap involving sampling 30 characters at random from among the 50 would have been more appropriate, though there is no way to know this in advance. For the purposes of this paper, we will ignore these correlations and assume that they cause no problems; in practice, they pose the most serious challenge to the use of bootstrap methods.

A similar problem can arise when multistate characters have been recoded into binary "factors" that are then treated as if they were independent two-state characters. These cannot be completely independent, as they would have to be if the bootstrap sampled them independently. Walter Fitch (pers. comm.) has suggested that this problem can be avoided by retaining a record of which binary factors are associated with which of the original characters, and having the bootstrap sample the original characters and keep all of the binary factors of a character together. Thus, if there were nine characters that had been expanded to 20 binary factors, we would construct the bootstrap sample by drawing nine times from the nine characters, and whenever a character was drawn we would take care to put all of its binary factors into the bootstrap sample.

Confidence Limits on Phylogenies

An interesting problem arises when we begin to consider how to construct confidence limits on the phylogenies. Each bootstrap sample is a data set that must

be analyzed to obtain an estimate of the phylogeny. We then have r phylogenies. Each of these is a complicated multivariate entity that has a tree topology and may also have branch lengths. Defining a confidence interval and summarizing it in a useable form is far from a simple matter.

In bootstrapping, confidence limits on a statistic are frequently constructed by the percentile method, which involves simply taking (for a 95% confidence interval) the empirical upper and lower 2.5% points of the distribution of bootstrap estimates of the statistic. Consider testing whether the probability of heads of a tossed coin exceeds 0.50. If we did not know about the binomial distribution and decided instead to use the bootstrap, a one-sided confidence interval on the probability of heads could be constructed by finding the empirical lower 5% point of the distribution of bootstrap estimates. The set of values less than 0.50 would therefore be rejected if values of the estimated probability of heads that small or smaller occurred less than 5% of the time among the bootstrap estimates.

The approach used here starts with the assumption that the systematist is primarily interested in whether some particular group is monophyletic. A rooted tree is a series of statements asserting monophyly of a series of nested or disjoint sets of species. Suppose that we are interested in a subset S of species and wish to know whether there is significant support in the data for the assertion that this group is monophyletic. We can reject the alternatives to the subset S if they occur in less than 5% of the bootstrap estimates.

We thus wish to search for all subsets S of species that occur on 95% or more of the bootstrap estimates. Each of these subsets may be considered to be supported (in the sense that its alternatives are rejected), although those confidence statements are not joint confidence statements: if two subsets are each supported at the 95% level, we might have as little as 90% confidence in the statement that

they are both present in the true tree. But at least they cannot be contradictory: each being present on at least 95% of the bootstrap estimated trees, they must co-occur on at least one of the trees and must thus be either nested or disjoint.

The same argument has been used by Margush and McMorris (1981) to define "majority rule" consensus trees. These are trees composed of all those subsets that appear in a majority of a collection of trees. By the argument just given, these subsets must define a tree, since no two of them can conflict. If we take the set of phylogenies that result from analyzing a series of bootstrap samples and make a majority-rule consensus tree, recording on it how often each subset appears, we will obtain a tree that can be used to define at a glance confidence sets for any rejection probability below 50%. The majority-rule consensus tree itself can be considered to be an overall bootstrap estimate of the phylogeny.

In cases where we are using a statistically well-founded method, such as maximum likelihood estimation, we would hope that the bootstrap method and the curvature of the likelihood surface would give similar indications of which parts of the phylogeny were well estimated and which not. Where the method of inferring phylogenies is one with undesirable statistical properties such as inconsistency, the bootstrap does not correct for these. For example, clustering by overall similarity makes an inconsistent estimate of the phylogeny if rates of evolution in different lineages differ by more than a certain amount. Parsimony methods are subject to the same problem, but require greater inequalities of evolutionary rate to be inconsistent. For an elementary discussion of these phenomena, see my recent review article (Felsenstein, 1983*b*). Bootstrapping provides us with a confidence interval within which is contained not the true phylogeny, but the phylogeny that would be estimated on repeated sampling of many characters from the underlying pool of characters. As such it may be misleading if the method used to infer phylogenies is inconsistent.

TABLE 1. Fossil horse data of Camin and Sokal (1965). The states of each character are in a linear series. -1, 0, 1, 2, . . . , with the ancestral state being 0. The data are also shown in binary recoded form in which the nine multistate characters have been recoded into 20 binary factors. The first line of that table indicates the correspondence between the original and recoded characters. Bootstrap sampling of characters should be done before any recoding into binary factors.

Name	Characters									Binary Factors			
										11112	22333	44566	77889
<i>Mesohippus</i>	0	0	0	0	0	0	0	0	0	00000	00000	00000	00000
<i>Hypohippus</i>	-1	3	3	0	0	0	0	0	1	00011	11111	00000	00001
<i>Archaeohippus</i>	1	0	0	0	0	0	0	1	1	10000	00000	00000	00101
<i>Parahippus</i>	1	1	1	1	0	0	0	0	1	10001	00100	10000	00001
<i>Merychippus</i>	2	2	2	2	1	1	1	2	1	11001	10110	11110	10111
<i>Merych. secundus</i>	2	2	2	2	1	-1	-1	2	1	11001	10110	11101	01111
<i>Nannippus</i>	2	2	1	2	1	1	1	2	1	11001	10100	11110	10111
<i>Neohipparion</i>	2	3	3	2	1	1	1	2	1	11001	11111	11110	10111
<i>Calippus</i>	2	2	1	2	1	-1	-1	2	1	11001	10100	11101	01111
<i>Pliohippus</i>	3	3	3	2	1	-1	-1	2	1	11101	11111	11101	01111

One difficulty in the interpretation of the result is that we may not have decided which subset of species interests us until after the bootstrap result is examined. This raises the "multiple tests" problem: if we have 20 statistical tests, on average one should show significance at the 95% level purely at random. There are ways of making simple corrections if the number of independent tests is known, but in this case the different tests (the different subsets that show up on the majority-rule consensus tree) are probably correlated, so that it is not easy to see how to compute the number of independent tests so as to correct for it. I have simply taken the 95% level as correct, as if we had chosen the test of interest a priori.

One might wonder whether the jackknife would be a viable alternative to the bootstrap. If we make a set of estimates by dropping one character at a time and then estimating the phylogeny, the resulting phylogenies will vary far less than the bootstrap estimates do. In the simple test case of sample averages estimating the mean of a normal distribution, it turns out that the jackknife estimates of the mean will have a variance only $n^2/(n-1)^3$ times as large as that of the corresponding bootstrap estimates (Efron and Gong, 1983). To make the variance among the jackknife estimates as large as that among bootstrap estimates,

one would have to engage in an extrapolation to make their variance larger. The difficulty in envisaging a procedure like this is that the space of possible phylogenies does not lend itself readily to extrapolation: once a branch length has shrunk to zero it is not immediately obvious what to do next. Unlike normal means, phylogenies do not live in a flat Euclidean space. One way to make the jackknife vary as much as the bootstrap would be to drop not one observation, but half the observations chosen at random. This possibility is worth exploring, but for the moment it is not obvious what advantage there would be to using the jackknife rather than the bootstrap.

Using Existing Computer Programs

The process of generating many bootstrap samples from a data set is a tedious one. One might think that it would require special programs to rewrite the data matrix, leaving out some characters and duplicating others. Fortunately, much of that work can be avoided by making use of differential character weights, which are allowed in most computer programs for inferring phylogenies, particularly programs using parsimony methods. These programs usually allow integer weights for the characters, weights that can be 0, 1, 2, A weight of zero means that the character is in effect

dropped from the analysis. A weight of w means that the character is counted as if present w times, so that each change of state in the character is counted as if it were w changes of state.

This automatically accomplishes the duplication and deletion of characters without the necessity of recopying the data matrix. Different bootstrap samples could be fed into the programs by doing computer runs with different weights. The weights are generated by starting with weights of zero for all characters. We then sample n characters at random with replacement (using a table of random numbers, for example). Each time a character is drawn, its weight is increased by one, so that in the end its weight counts the number of times it was sampled. Here are five of the weight vectors that were generated when bootstrap sampling was done on 20 characters:

```
21100212120121012010
01001510031020011211
01031211201012100121
12130421030000100101
20010100211211121211
```

To do bootstrap sampling, one would generate a vector of weights, run the phylogeny estimation program with those weights, generate another vector of weights, run the program with those, and so on. The process is fairly tedious, although with microcomputers it need not be expensive. In the fossil horse example below, I have used 50 bootstrap samples. This might strike a statistician as too few, but a systematist as too many. The more samples are taken, the more accurate an idea we will have of which groups are likely to be monophyletic. Even with a small number of bootstrap samples we will quickly get a feel for which parts of our estimate of the phylogeny are well supported and which not.

A computer program that carries out bootstrap sampling and computes the majority rule consensus tree is available for the case of discrete characters analyzed by the parsimony and compatibility methods. It is contained in the pro-

gram package PHYLIP, available free from me (see the Appendix below).

An Example

Table 1 shows the fossil horse data given by Camin and Sokal (1965 pp. 321–322) as a computational example. The full list of species and references for the original data are given by Camin and Sokal (1965). The data set has ten species and nine multistate characters. *Mesohippus* has been taken as the outgroup, as it was in Camin and Sokal's paper.

Figure 1 shows the results of running a branch-and-bound program that finds all most parsimonious trees according to the Wagner parsimony criterion. There are ten most parsimonious trees. The left tree in Figure 1 shows nine of them: the two empty circles with three descendants represent not trifurcations, but points at which the tree can be resolved into any of three bifurcating topologies. All nine possible combinations of these are in the list of most parsimonious trees. The tenth tree is the one shown at the right of Figure 1. All of these trees require 29 changes of character state.

If we were to take the variation among the most parsimonious trees as providing an adequate indication of the uncertainty in our estimate of the phylogeny, we would conclude that four monophyletic groups were defined, as these groups show up in all ten of the most parsimonious trees. They are:

- (*Pliohippus*, *Merychippus secundus*, *Calippus*)
- (*Nannippus*, *Neohipparion*, *Merychippus*)
- (*Pliohippus*, *Merychippus secundus*, *Calippus*, *Nannippus*, *Neohipparion*, *Merychippus*)
- (all but *Mesohippus* and *Archaeohippus*)

When we carry out bootstrap sampling of columns from the left-hand part of Table 1 and analyze 50 bootstrap replicates, we get the results shown in Figure 2. Next to each branch of the tree is shown the number of times that the bootstrap es-

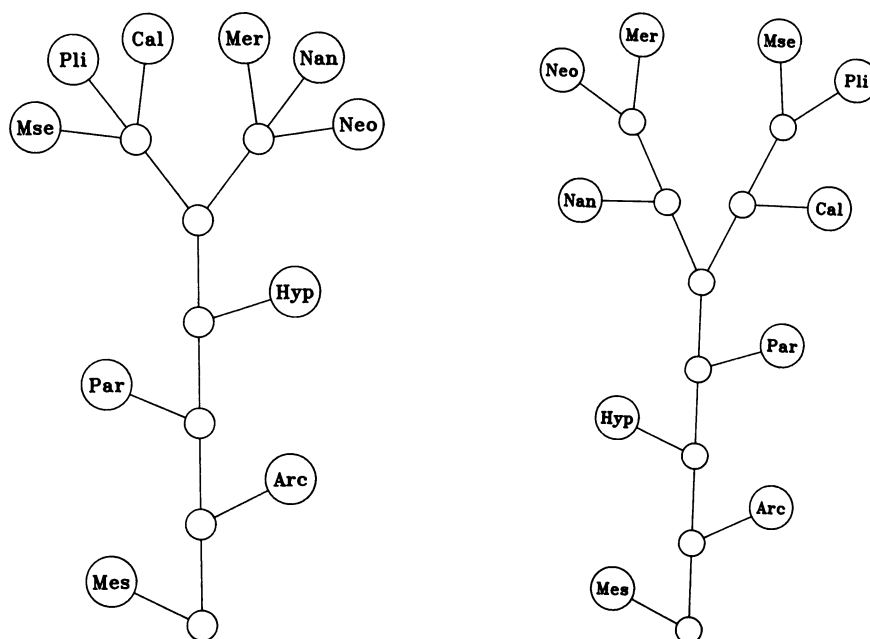


FIG. 1. All most parsimonious trees for the fossil horse data in Table 1 when phylogenies are evaluated by the Wagner parsimony criterion. There are ten most parsimonious trees in all. Nine of these can be generated by resolving each of the trifurcations in the left tree into all three possible bifurcations. The tenth is shown in the right tree. All abbreviations are first three letters of names in Table 1, except MSE = *Merychippus secundus*.

timate contained the corresponding monophyletic group. The tree shown is the majority-rule consensus tree. The consensus tree turns out in this case not to be one of the most parsimonious trees. All four of the monophyletic groups listed above occur on it, but only one of these (the six-species group consisting of *Pliohippus*, *Merychippus secundus*, *Calippus*, *Nannippus*, *Neohipparion*, and *Merychippus*) comes close to occurring 95% of the time in the bootstrap sampling (it occurs 47/50 or 94% of the time). The others only occur about two-thirds of the time. It is apparent that taking the set of most parsimonious trees as defining the confidence interval would result in far too narrow an interval.

The example given here has had the tree rooted by use of an outgroup. If we were using a method that produced an unrooted tree and had no outgroup in-

formation or other method of rooting the tree, we could still carry out bootstrap sampling and construct an unrooted majority-rule consensus tree. To do that, we would only need to note that each branch of one of the replicate bootstrap estimates divides the species into two groups, at least one of which would be monophyletic if we could root the tree. The unrooted majority-rule consensus tree is defined by finding those partitions that occur in a majority of the replicate trees. A simple way of doing this is to choose an arbitrary species as an outgroup, make a majority-rule consensus tree of the resulting rooted trees, and then present the result as an unrooted tree without indicating which species was the outgroup.

Perfectly Hennigian Data

Occasionally, though rather rarely, a data set will arise that has no internal

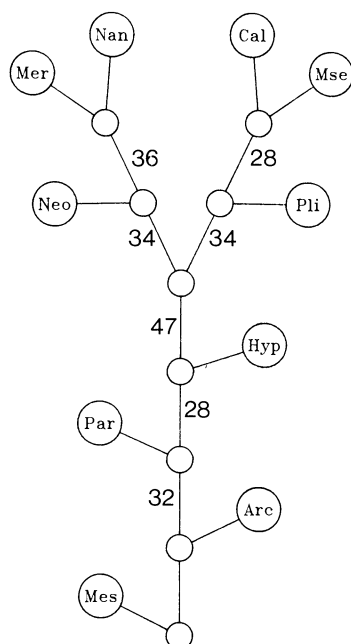


FIG. 2. Bootstrap estimate of the phylogeny for the data of Table 1 when phylogenies are evaluated by the Wagner parsimony criterion. Fifty bootstrap samples were analyzed. The groups shown in this tree are those that occurred in a majority of the resulting trees, plus the most frequently occurring groups that were compatible with these. Next to each branch is shown the number of times that the monophyletic group it defines occurred. Further explanation is given in the text.

conflict, all characters being perfectly compatible. This sort of data, which is that envisioned by Hennig when he suggested using derived character states to define monophyletic groups, allows us to avoid entirely the bootstrap sampling process. The argument is quite simple. Suppose that we want to construct a 95% confidence interval by bootstrap sampling. Suppose that c characters out of n define the same monophyletic group. That group will show up in the bootstrap estimate if any one of those c characters is drawn in the sampling of n characters. The monophyletic group will be part of the 95% confidence interval if and only if the probability of omitting all c of the

characters is less than 0.05. This is easily computed, given n and c .

The probability of leaving out all c characters in drawing n characters without replacement is $(1 - c/n)^n$. The value of c that is necessary to make this less than 0.05 is the same for all relevant values of n : it is $c = 3$. We can thus conclude that, if the data are perfectly Hennigian, three characters are enough for the bootstrap to indicate significant support for a monophyletic group at the 95% level. Any group supported by fewer characters will not be in the bootstrap confidence interval. Of course, we are assuming that the evolutionary processes and the inclusion of the characters by the systematist are independent across characters.

Although three characters are enough to guarantee inclusion of a group, if the data are perfectly Hennigian, one will never encounter any character that contradicts the group. Sometimes we have great confidence that our characters are "clean" ones, that reversals and parallelisms would be so rarely seen that we can have confidence in a group even if it is supported by only one character. The present "rule of three" would then seem to be a conservative one.

It may be doubted that the rule is really always conservative. I have recently studied, by exact enumeration methods, the problem of placing confidence limits on phylogenies using parsimony methods when there are only three species and an evolutionary process for which an evolutionary clock may be assumed (Felsenstein, 1985). It turns out that in the worst case, when the characters are equally likely to resolve a trifurcation in any of the three possible ways, if we have three characters all of which support the same resolution, this is not statistically significant at the 95% level. Four characters would be. (I am indebted to Alan Templeton for pointing out the connection between the two calculations.)

In many cases, strong conclusions have been drawn from the existence of groups defined by as little as one character. The great advantage of the present approach

is that it provides a practical method, albeit a flawed one, for assessing the uncertainty inherent in such conclusions. I suspect that the levels of uncertainty found in practice will be so great as to give pause to all but the firmest exponents of nonstatistical hypothetico-deductive approaches to inferring phylogenies.

ACKNOWLEDGMENTS

I am grateful to Kent Fiala of the Department of Ecology and Evolution, State University of New York at Stony Brook, for providing me with the fossil horse data of Camin and Sokal in recoded form. I wish to thank Walter Fitch, Monty Slatkin, Alan Templeton, Bill Engels, Ruth Shaw, and an anonymous statistical reviewer for suggestions for improvement of the manuscript. This work was supported by task agreement number DE-AT06-76EV71005 of contract number DE-AM06-76RL02225 between the U.S. Department of Energy and the University of Washington.

LITERATURE CITED

- CAMIN, J. H., AND R. R. SOKAL. 1965. A method for deducing branching sequences in phylogeny. *Evolution* 19:311-326.
- CAVENDER, J. A. 1978. Taxonomy with confidence. *Math. Biosci.* 40:271-280 (Erratum: Vol. 44, p. 308, 1979).
- . 1981. Tests of phylogenetic hypotheses under generalized models. *Math. Biosci.* 54:217-229.
- DIACONIS, P., AND B. EFRON. 1983. Computer-intensive methods in statistics. *Sci. Amer.* 249: 116-130.
- EFRON, B. 1979. Bootstrap methods: Another look at the jackknife. *Ann. Statist.* 7:1-26.

- . 1982. The jackknife, the bootstrap, and other resampling plans. CBMS-NSF Regional Conference Series in Applied Mathematics No. 38. Society for Industrial and Applied Mathematics. Philadelphia, PA.
- EFRON, B., AND G. GONG. 1983. A leisurely look at the bootstrap, the jackknife, and cross-validation. *Amer. Statist.* 37:36-48.
- FELSENSTEIN, J. 1983a. Statistical inference of phylogenies. *J. Roy. Statist. Soc. A* 146:246-272.
- . 1983b. Parsimony in systematics: Biological and statistical issues. *Ann. Rev. Ecol. Syst.* 14:313-333.
- . 1985. Confidence limits on phylogenies with a molecular clock. *Systematic Zoology* 34: 152-161.
- MARGUSH, T., AND F. R. MCMORRIS. 1981. Consensus n-trees. *Bull. Mathemat. Biol.* 43:239-244.
- TEMPLETON, A. R. 1983. Phylogenetic inference from restriction endonuclease cleavage site maps with particular reference to the evolution of humans and the apes. *Evolution* 37:221-224.

Corresponding Editor: W. R. Engels

APPENDIX

Availability of the PHYLIP Program Package

PHYLIP, the Phylogeny Inference Package, is a free package of computer programs, written in Pascal, for inferring phylogenies. It includes parsimony methods, compatibility methods, distance matrix methods, and maximum likelihood methods. The Pascal source code is provided (compiled object code is not). PHYLIP will be written in a standard format on a magnetic tape provided by the recipient. It will also be provided on 5¼-inch diskettes if 6 double density diskettes are sent. A variety of soft-sectored MSDOS, CP/M-80, and CP/M-86 formats can be written; double-sided, hard-sectored, or 3.5-inch formats cannot, nor can any Apple formats. For information on formats supported and restrictions on countries to which distribution and support are available, please write the author.