

The asymptotic behavior of bootstrap support values in molecular phylogenetics

JUN HUANG^{1,2}, YUTING LIU^{1,*}, TIANQI ZHU^{3,*} AND ZIHENG YANG^{2,*}

¹*Department of Mathematics, Beijing Jiaotong University, Beijing, 100044, China*

²*Department of Genetics, Evolution and Environment, University College London, Gower Street, London WC1E 6BT, UK.*

³*Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Beijing 100190, China.*

* *Correspondence to be sent to: Yuting Liu, Department of Mathematics, Beijing Jiaotong University, Beijing, 100044, China; E-mail: ytliu@bjtu.edu.cn*

Tianqi Zhu, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Beijing 100190, China; E-mail: zhutq@amss.ac.cn

Ziheng Yang, Department of Genetics, Evolution and Environment, University College London, Gower Street, London WC1E 6BT, UK; E-mail: z.yang@ucl.ac.uk (orcid: 0000-0003-3351-7981)

ABSTRACT

1 The phylogenetic bootstrap is the most commonly used method for assessing statistical
2 confidence in estimated phylogenies by non-Bayesian methods such as maximum parsimony and
3 maximum likelihood (ML). It is observed that bootstrap support tends to be high in large
4 genomic datasets whether or not the inferred trees and clades are correct. Here we study the
5 asymptotic behavior of bootstrap support for the ML tree in large datasets when the competing
6 phylogenetic trees are equally right or equally wrong. We consider phylogenetic reconstruction as
7 a problem of statistical model selection when the compared models are nonnested and
8 misspecified. The bootstrap is found to have qualitatively different dynamics from Bayesian
9 inference, and does not exhibit the polarized behavior of posterior model probabilities, consistent
10 with the empirical observation that the bootstrap is more conservative than Bayesian
11 probabilities. Nevertheless bootstrap support similarly shows fluctuations among large datasets,
12 with no convergence to a point value, when the compared models are equally right or equally
13 wrong. Thus in large datasets strong support for wrong trees or models is likely to occur. Our
14 analysis provides a partial explanation for the high bootstrap support values for incorrect clades
15 observed in empirical data analysis.

16 *Key words:* Bootstrap, model selection, star-tree paradox, support value

INTRODUCTION

17
18 Recently Yang and Zhu (2018) characterized the asymptotic behaviors of Bayesian model
19 selection in large datasets. When two models are both right or are equally wrong and indistinct,
20 the posterior model probability varies among datasets according to a statistical distribution such
21 as $\mathbb{U}(0, 1)$, whereas one might expect it to converge to the point value $\frac{1}{2}$. Even more disturbingly,
22 when the two models are equally wrong and distinct, the posterior model probability approaches
23 $\sim 100\%$ in some datasets and 0% in others. This polarized behavior may be a major reason for
24 the observation that in Bayesian analysis of large phylogenetic datasets, posterior probabilities
25 for trees or clades are most often close to 100% , whether or not the relationships are correct.

26 For non-Bayesian methods including maximum parsimony (Fitch, 1971), neighbor
27 joining (Saitou and Nei, 1987), and maximum likelihood (ML, Felsenstein, 1981), confidence for
28 inferred trees or clades is most often assessed using Felsenstein’s (1985) phylogenetic bootstrap.
29 An interesting question is whether bootstrap exhibits similar behaviors as the posterior model
30 probabilities. In modern phylogenomic studies both posterior probabilities and bootstrap support
31 values tend to be $\sim 100\%$, whether or not the clades or trees are correct. Such results lead to
32 widespread mistrust for bootstrap support values in large datasets. For example, Chan *et al.*
33 (2020) wrote that “high bootstrap support did not necessarily reflect congruence or support for
34 the correct topology. This study reiterates findings of some previous studies, which demonstrated
35 that traditional bootstrap values can produce positively misleading measures of support in large
36 phylogenomic datasets.”

37 Bootstrap was originally developed by Efron (1979) to calculate the standard error for a
38 parameter, by resampling the original data and studying the variation among the bootstrap
39 resample datasets. It has since been used to conduct all sorts of analyses in Frequentist statistics,
40 such as correction for bias, calculation of standard errors, construction of confidence intervals,
41 and performing significance tests (Efron and Tibshirani, 1993; Davison and Hinkley, 1997). In
42 phylogenetics, bootstrap was introduced by Felsenstein (1985) to assess the confidence in
43 estimated phylogenetic trees. Although it follows the same operational procedure of resampling
44 data points from the observed dataset, bootstrap in phylogenetics differs from its use in bias
45 correction or in confidence-interval construction, in that a statistical interpretation has been
46 illusory despite numerous efforts (Zharkikh and Li, 1992; Hillis and Bull, 1993; Felsenstein and
47 Kishino, 1993; Berry and Gascuel, 1996; Efron *et al.*, 1996; Holmes, 2003; Susko, 2009).
48 Modifications to the procedure have also been made, including the complete-and-partial
49 bootstrap (Zharkikh and Li, 1995), correction for first-order biases (Susko, 2010), or adjustment
50 for short branches (Lemoine *et al.*, 2018). These correct for the perceived bias in the procedure or
51 to make it agree better with standard ideas of confidence levels and hypothesis testing.

52 Its interpretation aside, phylogenetic bootstrap is the most widely used procedure for
53 assessing the confidence in estimated phylogenies by non-Bayesian methods. Felsenstein’s 1985
54 paper is a citation classic in all sciences. For Bayesian methods, the posterior probability for the
55 inferred tree provides a natural measure of uncertainty (Rannala and Yang, 1996), and bootstrap
56 is in theory not needed in Bayesian inference. However, the sensitivity of Bayesian model choice
57 to the prior (O’Hagan and Forster, 2004) and the polarized behavior of Bayesian model selection
58 under model misspecification (Yang and Zhu, 2018) have prompted the application of bootstrap
59 in Bayesian model selection as well, leading to methods such as Bayesian bagging (Rubin, 1981;
60 Weng, 1989; Huggins and Miller, 2020). It is important to study the asymptotic behavior of
61 phylogenetic bootstrap. Earlier simulation studies suggest that the phylogenetic bootstrap may be
62 conservative, and that 70% (instead of 95%) means strong support (e.g., Hillis and Bull, 1993). It
63 has been noted that bootstrap support is numerically less extreme than posterior model
64 probabilities (e.g., Huelsenbeck and Rannala, 2004; Yang and Rannala, 2005).

65 In this paper we explore the asymptotic behavior of phylogenetic bootstrap when the data
66 size increases. We consider phylogenetic reconstruction as a statistical model selection problem,
67 and treat phylogenetic trees as nonnested statistical models (rather than different values of a
68 parameter in a well-specified model). We present an asymptotic theory for bootstrap model
69 probability under different scenarios in the Appendix, and in the main paper illustrate the theory
70 using canonical problems that are analytically tractable. We discuss phylogenetic reconstruction
71 problems in the case of three or four taxa to illustrate the general theory.

72 SUMMARY OF ANALYTICAL RESULTS

73 Following Felsenstein and Kishino (1993) and Efron *et al.* (1996), we consider bootstrap as a
74 general approach to assessing the confidence in the selected model in a model-selection problem.

Bootstrap in model selection

75

76 The data are an independently and identically distributed (i.i.d.) sample of size n , $x =$
 77 $\{x_1, \dots, x_n\}$, from the true data-generating model $g(X)$. We compare K models, H_j , $j = 1, \dots, K$.
 78 Model H_j specifies the density $f_j(X|\theta_j)$ with parameters θ_j . Let $\hat{\theta}_j$ be the MLE of θ_j under
 79 model H_j given data x . When $n \rightarrow \infty$, $\hat{\theta}_j \rightarrow \theta_{j*}$, where θ_{j*} minimizes the Kullback-Leibler (K-L)
 80 divergence from model H_j to the true model,

$$D_j = \int g(X) \log \frac{g(X)}{f_j(X|\theta_{j*})} dX. \quad (1)$$

81 If H_j is correct, θ_{j*} will be the *true parameter values*, with $D_j = 0$. Otherwise if H_j is wrong, θ_{j*}
 82 will be the *best-fitting* or *pseudo-true parameter values*, with $D_j > 0$. In this paper we focus on
 83 the case where all K models have the same K-L divergence to the true model. Two models f_1 and
 84 f_2 are said to be equally right if $D_1 = D_2 = 0$, and equally wrong if $D_1 = D_2 > 0$. If two models
 85 are unidentifiable at their pseudo-true parameter values, that is, if

$$f_1(X|\theta_{1*}) = f_2(X|\theta_{2*}) \quad \text{for almost every } X, \quad (2)$$

86 they are said to be indistinct. This can occur when both models are right (with $D_1 = D_2 = 0$) or
 87 when both are wrong (with $D_1 = D_2 > 0$). Otherwise if equation (2) does not hold for some X of
 88 nonzero measure, the models are said to be distinct. This can occur only if both models are wrong
 89 (with $D_1 = D_2 > 0$).

90 The model selected by ML is the one that achieves the greatest log likelihood, $\ell_j(\hat{\theta}_j) =$
 91 $\log f_j(x|\hat{\theta}_j)$. To assess the confidence on the selected model, we calculate the bootstrap
 92 probability. Let $x_b^* = \{x_{b1}^*, \dots, x_{bn}^*\}$ be a bootstrap sample, formed by resampling with
 93 replacement n times from the original data x . Let $\hat{\theta}_b^*$ be the MLE from a bootstrap sample x^* .
 94 Here we follow the convention of using the superscript $*$ to indicate a bootstrap sample, and the
 95 subscript $*$ for the true or pseudo-true parameter values. We assume that θ_{j*} , $\hat{\theta}_j$, and $\hat{\theta}_j^*$ are inner
 96 points in the parameter space. The proportion of bootstrap replicates in which model j is the
 97 optimal model is the bootstrap probability or bootstrap support P_j for model j . For example, in
 98 the case of two models, the bootstrap probability for model H_1 is

$$P_1(x) = \mathbb{P} \{ \log f_1(x^*|\hat{\theta}_1^*) > \log f_2(x^*|\hat{\theta}_2^*) | x \} \approx \frac{1}{B} \sum_b \mathbb{I}_{\ell_1(\hat{\theta}_1^*) > \ell_2(\hat{\theta}_2^*)}, \quad (3)$$

99 where $\ell_j(\hat{\theta}_j^*) = \log f_j(x^*|\hat{\theta}_j^*)$ is the log likelihood value for model j , calculated at the MLE ($\hat{\theta}_j^*$)
 100 and where the indicator function \mathbb{I}_A is 1 if A is true or 0 otherwise. Note that P_1 is a function of x
 101 and is a random variable. We are interested in the asymptotic distribution of P_1 when x varies.

102 In phylogenetics, the models under comparison are the tree topologies for the given set of
 103 species, while each data point corresponds to one site or one column in the alignment. While the
 104 bootstrap is applicable as long as the inference method is statistically consistent (Felsenstein,
 105 1985), we focus on ML in this paper. In phylogenetics, bootstrap is commonly used to attach
 106 support values for clades or splits on the phylogeny, calculated as the proportion of bootstrap
 107 trees that contain the splits. Here we focus on the bootstrap probability for the whole model. In
 108 the case of simple trees with three or four species with only one internal branch, the two
 109 measures are equivalent. We assume that the number of bootstrap replicates B is large so that the
 110 sampling errors due to limited number of bootstrap replicates is negligible.

111 *The asymptotic behavior of bootstrap model selection under different scenarios*

112 We develop an asymptotic theory of bootstrap model selection in the Appendix. In general, when
 113 equally right or equally wrong models are compared, bootstrap model probabilities have a
 114 non-degenerate distribution. In the case of two equally wrong and distinct models, the bootstrap
 115 model probability P_1 has the distribution $\mathbb{U}(0, 1)$.

116 The case of two equally wrong and distinct models with no parameters provides valuable
 117 insights into the differences between bootstrap and Bayesian methods. The log-likelihood ratio
 118 between the two models is

$$\Delta = \log \frac{f_1(x)}{f_2(x)}, \quad \Delta^* = \log \frac{f_1(x^*)}{f_2(x^*)}, \quad (4)$$

119 for the original data x and the bootstrap resample data x^* , respectively. Each of these is a sum of n
 120 i.i.d. terms. Thus $\mathbb{E}(\Delta) = n \mathbb{E} \log f_1(X) - n \mathbb{E} \log f_2(X) = n(D_2 - D_1) = 0$ (Equation 1). Let

$$\sigma^2 = \mathbb{V} \left\{ \log \frac{f_1(X)}{f_2(X)} \right\} = \int g(X) \left[\log \frac{f_1(X)}{f_2(X)} \right]^2 dX. \quad (5)$$

121 When $n \rightarrow \infty$, $\Delta \sim \mathbb{N}(0, n\sigma^2)$ and $\Delta^* | x \sim \mathbb{N}(\Delta, n\sigma^2)$, according to the central limit theorem. Thus

$$P_1 = \mathbb{P}\{\Delta^* > 0 | x\} = \Phi\left(\frac{\Delta}{\sqrt{n\sigma}}\right) \rightarrow \mathbb{U}(0, 1), \quad (6)$$

122 where Φ is the cumulative distribution function (CDF) for $\mathbb{N}(0, 1)$.

123 In Bayesian comparison of two equally wrong models with no parameters, Δ is the log
 124 Bayes factor. With equal prior probabilities ($\frac{1}{2}$ for each model), this is related to the posterior
 125 model probability through $\Delta = \log \frac{P_1}{1-P_1}$ or $P_1 = \frac{e^\Delta}{e^\Delta + 1}$. As Δ behaves like a random walk when n
 126 increases, it is nearly impossible for Δ to be in a small interval around 0, say, $-5 < \Delta < 5$ which
 127 corresponds to $0.007 < P_1 < 0.993$. In other words, for large n , the posterior probability will be 0
 128 in half of the datasets and 1 in the other half. This polarized behavior also occurs when the
 129 compared models, equally wrong and distinct, have parameters as the Bayes factor is dominated
 130 by the random-walk term (Yang and Zhu, 2018). The analysis here suggests that bootstrap
 131 probability has a qualitatively different behavior, as it contrasts Δ^* for the bootstrap sample with
 132 Δ for the original data.

133 1. ILLUSTRATIVE EXAMPLES

134 We present several simple examples to illustrate the asymptotic behavior of bootstrap model
 135 probability under different scenarios when the data size $n \rightarrow \infty$. In the first two examples two
 136 models are equally wrong and distinct, and the bootstrap probability P_1 varies among datasets
 137 like a random number, $P_1 \sim \mathbb{U}(0, 1)$ (Equation 6).

138 **Problem 1 fair-coin paradox, with equally wrong models and no parameters.**

139 Suppose a coin is fair with the true probability of heads to be $p = 0.5$, and we flip the coin n
 140 times to compare two models $H_1 : p = 0.4$ and $H_2 : p = 0.6$. The dataset is $x = \{x_1, \dots, x_n\}$,
 141 where x_i takes the value 1 for heads and 0 for tails, and has the bernoulli distribution. The data
 142 can be summarized as the proportion of heads in n tosses, \bar{x} , which is approximately normal
 143 $\mathbb{N}(\frac{1}{2}, \frac{1}{4n})$. H_1 is favored if $\bar{x} < \frac{1}{2}$, and this happens in half of the datasets.

144 Given x , the bootstrap sample $x_b^* = \{x_{b1}^*, \dots, x_{bn}^*\}$, where x_{bi}^* is bernoulli with probability
 145 \bar{x} , can be summarized as the bootstrap sample mean \bar{x}^* , which is approximately normal, with
 146 $\bar{x}^* | x \sim \mathbb{N}(\bar{x}, \frac{\bar{x}(1-\bar{x})}{n}) \approx \mathbb{N}(\bar{x}, \frac{1}{4n})$. The bootstrap sample x_b^* favors model H_1 if and only if $\bar{x}^* < 1/2$.
 147 Thus

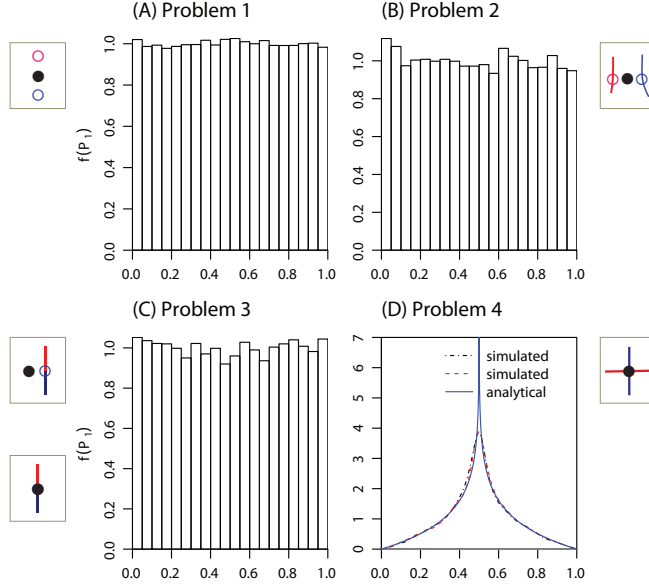


Fig. 1. Histogram/density of bootstrap model probability P_1 in comparisons of two models. (A) Problem 1 (the fair-coin paradox) in which a fair coin (with $p = 0.5$) is tossed n times to compare two equally wrong and distinct models: $p = 0.4$ and $p = 0.6$. (B) Problem 2 in which the true model is $\mathbb{N}(0, 1)$ while the two compared models, $\mathbb{N}(\mu, 1/\tau_1)$ and $\mathbb{N}(\mu, 1/\tau_2)$ with $\tau_1 < 1 < \tau_2$, are equally wrong and distinct. (C) Problem 3 (the fair-balance paradox) where the true model is $\mathbb{N}(0, 1)$ and the two compared models, $\mathbb{N}(\mu, 1/\tau)$, $\mu < 0$ versus $\mathbb{N}(\mu, 1/\tau)$, $\mu > 0$, are equally right (if $\tau = 1$) or equally wrong and indistinct (if $\tau \neq 1$). (D) Problem 4 (equally right models). The true model is $\mathbb{N}(0, 1)$ while the two compared models, $\mathbb{N}(\mu, 1)$ versus $\mathbb{N}(0, 1/\tau)$, are both right. Black dashed line is for the expensive simulation generating x and x^* , the red dashed line is for simulation generating \bar{x} and s^2 , while blue solid line is for the analytical approximation by Equation 14. The insets characterize the problems, with the true models represented as filled circles and the pseudo-true parameter values as empty circles, while the lines represent the parameter space for each model. The settings are $n = 10^5$, $B = 3 \times 10^4$, and $R = 10^5$ for problem 1, $n = 10^4$, $B = 3 \times 10^4$, and $R = 10^4$ for problem 2, $n = 10^4$, $B = 3 \times 10^4$, and $R = 10^4$ for problem 3, and $n = 10^4$, $B = 10^3$, and $R = 10^4$ for problem 4.

$$P_1 = \mathbb{P}\{\bar{x}^* < \frac{1}{2}|x\} \approx \Phi\left(\frac{1/2 - \bar{x}}{\sqrt{1/(4n)}}\right) \rightarrow \mathbb{U}(0, 1), \quad \text{as } n \rightarrow \infty. \quad (7)$$

Thus P_1 varies like a random number among datasets (Fig. 1A). Alternatively we have $\Delta = \ell_1 - \ell_2 = 2n(\bar{x} - \frac{1}{2}) \log \frac{0.4}{0.6} \sim \mathbb{N}(0, n\sigma^2)$ and $\Delta^*|x \sim \mathbb{N}(\Delta, n\sigma^2)$, with $\sigma = \log \frac{0.4}{0.6}$, so that Equation 6 gives $P_1 \sim \mathbb{N}(0, 1)$.

Problem 2 Normal distribution, equally wrong and distinct models with free parameters. Suppose the true model is $\mathbb{N}(0, 1)$ and we consider $H_1 : \mathbb{N}(\mu, 1/\tau_1)$ and $H_2 : \mathbb{N}(\mu, 1/\tau_2)$, where μ is a free parameter while the precisions τ_1 and τ_2 are given with $\log(\tau_2/\tau_1) = \tau_2 - \tau_1$ so that the two models are equally wrong ($D_1 = D_2 > 0$) (Yang and Zhu, 2018). We use $\tau_1 = 0.25$ and $\tau_2 = 2.58666$. Under each model, the pseudo-true parameter value is $\mu_* = 0$ and H_1 and H_2 are two equally wrong and distinct models. Note that H_1 is over-dispersed and H_2 is under-dispersed. Under the model $\mathbb{N}(\mu, 1/\tau)$ with known τ , the log likelihood is

$$\ell = -\frac{n}{2} \log(2\pi) + \frac{n}{2} \log \tau - \frac{\tau}{2} \sum_{i=1}^n (x_i - \mu)^2, \quad (8)$$

with $\hat{\mu} = \bar{x}$. Thus $\ell_1 > \ell_2$ if and only if $(\tau_2 - \tau_1) \sum_{i=1}^n (x_i - \bar{x})^2 > n \log(\tau_2/\tau_1)$ or if and only if $s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 > 1$. We have $ns^2 \sim \chi_{n-1}^2 \approx \mathbb{N}(n-1, 2(n-1))$.

Given x , the bootstrap sample x^* favors H_1 if the sample variance $s^{2*} > 1$. We have $ns^{2*}/s^2|x \sim \chi_{n-1}^2 \approx \mathbb{N}(n-1, 2(n-1))$. For large n , re-sampling from the empirical distribution represented by the observed data x is approximately equivalent to sampling from the continuous distribution $\mathbb{N}(\bar{x}, s^2)$. Thus

$$P_1 = \mathbb{P}\{s^{2*} > 1|x\} \approx \Phi\left(\frac{((n-1)/n)s^2 - 1}{\sqrt{2(n-1)/n^2}}\right) \rightarrow \mathbb{U}(0, 1), \quad \text{as } n \rightarrow \infty. \quad (9)$$

This is confirmed in Fig. 1B.

Alternatively we have $\Delta = \ell_1 - \ell_2 = \frac{n}{2}(\tau_2 - \tau_1)(s^2 - 1) \sim \mathbb{N}(0, n\sigma^2)$ and $\Delta^*|\Delta \sim \mathbb{N}(\Delta, n\sigma^2)$, with $\sigma = \frac{1}{\sqrt{2}}(\tau_2 - \tau_1)$, so that $P_1 = \mathbb{P}\{\Delta^* > 0|x\} \sim \mathbb{N}(0, 1)$.

If the two compared models are both right (with $D_1 = D_2 = 0$) or are equally wrong and indistinct (with $D_1 = D_2 > 0$), then P_1 varies among datasets according to a nondegenerate distribution, which may and may not be $\mathbb{U}(0, 1)$, as illustrated in the next two examples.

Problem 3 (fair-balance paradox with two equally right or equally wrong and indistinct models). The true model is $\mathbb{N}(0, 1)$ and the two compared models are $\mathbb{N}(\mu, 1/\tau)$, $\mu < 0$ and $\mathbb{N}(\mu, 1/\tau)$, $\mu > 0$, with τ given. If $\tau = 1$, the two models are equally right. If $\tau \neq 1$, the two models are equally wrong (because of the assumed incorrect variance) and indistinct (because the pseudo-true parameter value $\mu_* = 0$ under each model). Model 1 is favored if and only if the sample mean $\bar{x} < 0$. As $\bar{x} \sim \mathbb{N}(0, 1/n)$ and $\bar{x}^*|x \sim \mathbb{N}(\bar{x}, 1/n)$, we have

$$P_1 = \mathbb{P}\{\bar{x}^* < 0|x\} = \Phi(-\sqrt{n}\bar{x}) \rightarrow \mathbb{U}(0, 1), \quad \text{as } n \rightarrow \infty. \quad (10)$$

This is confirmed in Fig. 1C.

Problem 4 Normal-distribution example with an infinite spike at $\frac{1}{2}$ in the P_1 distribution. The true model is $\mathbb{N}(0, 1)$ and the two compared models are $\mathbb{N}(\mu, 1)$ and $\mathbb{N}(0, 1/\tau)$. In H_1 , $\mu_* = 0$ while in H_2 , $\tau_* = 1$, so the two models are equally right. The data x may be summarized as the sample mean \bar{x} and sample variance $s^2 = \frac{1}{n}\sum_i(x_i - \bar{x})^2$. The MLE of the parameter is $\hat{\mu} = \bar{x}$ under H_1 and $\hat{\tau} = n/\sum_i x_i^2 = 1/(s^2 + \bar{x}^2)$ under H_2 . The log-likelihood values are

$$\begin{aligned} \ell_1(\hat{\mu}) &= -\frac{1}{2}\sum(x_i - \bar{x})^2 = -\frac{1}{2}ns^2, \\ \ell_2(\hat{\tau}) &= -\frac{n}{2}\log\left(\frac{1}{n}\sum x_i^2\right) - \frac{n}{2} = -\frac{n}{2}\log(s^2 + \bar{x}^2) - \frac{n}{2} \end{aligned} \quad (11)$$

Thus $\ell_1 > \ell_2$ if and only if

$$\bar{x}^2 > e^{s^2-1} - s^2 \approx 1 + (s^2 - 1) + \frac{1}{2}(s^2 - 1)^2 - s^2 = \frac{1}{2}(s^2 - 1)^2, \quad (12)$$

or if and only if

$$|\bar{x}| > \frac{1}{\sqrt{2}}|s^2 - 1|. \quad (13)$$

A large deviation of \bar{x} from 0 supports H_1 , whereas a large deviation of s^2 from 1 favors H_2 . Also $\bar{x} \sim \mathbb{N}(0, \frac{1}{n})$ and $s^2 \sim \frac{1}{n}\chi_{n-1}^2 \approx \mathbb{N}(\frac{n-1}{n}, \frac{2(n-1)}{n^2})$ or $\frac{1}{\sqrt{2}}(s^2 - 1) \sim \mathbb{N}(0, \frac{1}{n})$, and \bar{x} and s^2 are independent. Thus Equation 13 holds and H_1 is the selected model in half of the datasets.

Given x , we have $\bar{x}^*|x \sim \mathbb{N}(\bar{x}, s^2/n) \approx \mathbb{N}(\bar{x}, \frac{1}{n})$ and $\frac{1}{\sqrt{2}}(s^{2*} - 1)|x \sim \mathbb{N}(\frac{1}{\sqrt{2}}(s^2 - 1), \frac{1}{n})$ and \bar{x}^* and s^{2*} are conditionally independent. Let $z_1 = \sqrt{n}\bar{x}$ and $z_2 = \sqrt{\frac{n}{2}}(s^2 - 1)$, with z_1 and z_2 from $\mathbb{N}(0, 1)$. Let $z_1^* = \sqrt{n}\bar{x}^*$ and $z_2^* = \sqrt{\frac{n}{2}}(s^{2*} - 1)$, with $z_1^*|x \sim \mathbb{N}(z_1, 1)$ and $z_2^*|x \sim \mathbb{N}(z_2, 1)$ to be conditionally i.i.d. Then

$$P_1 = \mathbb{P}\{|z_1^*| > |z_2^*||x\}. \quad (14)$$

192 This problem is analyzed in the SI text, available on Dryad at
 193 <https://doi.org/10.5061/dryad.7m0cfxprw>. The limiting distribution of P_1 when $n \rightarrow \infty$ is

$$f(P_1) = -\log|2P_1 - 1|. \quad (15)$$

194 The density is symmetrical around $\frac{1}{2}$, is 0 at 0 and 1, and has an infinite spike at $\frac{1}{2}$, with
 195 the mean $\frac{1}{2}$ and variance $\frac{1}{36}$. This is confirmed by simulation in Fig. 1D. The simulation is done
 196 in two ways. In the first, data x is sampled from $\mathbb{N}(0, 1)$, and given x bootstrap samples x_b^* are
 197 generated, with \bar{x}^* and s^{2*} calculated to apply Equation 14. In the second approach, $\bar{x} \sim \mathbb{N}(0, 1/n)$
 198 and $ns^2 \sim \chi_{n-1}^2$ are sampled, and then $\bar{x}^* \sim \mathbb{N}(\bar{x}, s^2/n)$ and $ns^{2*}/s^2 \sim \chi_{n-1}^2$ are generated to select
 199 the model for the bootstrap sample using Equation 14. Both approaches produce the same results
 200 as Equation 15.

201 **Problem 5 multivariate normal-distribution example.** The true model is the
 202 $(K - 1)$ -variate normal distribution $\mathbb{N}(\mu, \Sigma)$, with mean vector $\mu = (\mu_1, \dots, \mu_{K-1})$ where $\mu_1 =$
 203 $\dots = \mu_{K-1} = 0$ and variance matrix Σ which has 1 on the diagonal and $-1/(K - 1)$ on the
 204 off-diagonal. The data are an i.i.d. sample of size n , $x = \{x_{ij}\}$, $i = 1, \dots, n$; $j = 1, \dots, K - 1$. Also
 205 let $x_{iK} = -(x_{i1} + \dots + x_{i,K-1})$ and $\mu_K = -(\mu_1 + \dots + \mu_{K-1})$. We use the data to compare K
 206 models. Model H_j , $j = 1, \dots, K$, assumes $\mu_j > \mu_k$ for any $k \neq j$. The model has $K - 1$ free
 207 parameters: μ_1, \dots, μ_K with the constraint $\mu_1 + \dots + \mu_K = 0$. The variance is assumed to be
 208 known, $c\Sigma$. The models are equally right if $c = 1$ and equally wrong if $c \neq 1$. An alternative
 209 formulation of the problem is to have only one parameter in model H_j : $\mu_j > \mu_k$ with
 210 $\mu_k = -\mu_j/(K - 1)$ for all $k \neq j$.

211 Let $\bar{x} = \{\bar{x}_j\}$ and $\bar{x}^* = \{\bar{x}_j^*\}$, with

$$\bar{x}_j = \frac{1}{n} \sum_i x_{ij}, \quad \bar{x}_j^* = \frac{1}{n} \sum_i x_{ij}^*, \quad j = 1, \dots, K, \quad (16)$$

212 be the sample means from dataset x and from bootstrap sample x^* , respectively. Then
 213 $\bar{x} \sim \mathbb{N}(\mu, \frac{1}{n}\Sigma)$ and approximately $\bar{x}^*|x \sim \mathbb{N}(\bar{x}, \frac{1}{n}\Sigma)$. Without the constraint under each model H_j :
 214 $\mu_j > \mu_k$, the MLEs of μ are the sample means. With the constraint, H_j is the selected model if \bar{x}_j
 215 is the greatest among $\bar{x}_1, \dots, \bar{x}_K$. The bootstrap probability for model H_1 given data x is

$$P_1 = \mathbb{P}(\bar{x}_1^* > \bar{x}_2^*, \dots, \bar{x}_1^* > \bar{x}_K^* | x). \quad (17)$$

216 Now for any $j \neq k$,

$$\sigma_{jk}^2 = \mathbb{V}(\bar{x}_j - \bar{x}_k) = \frac{2}{n} - 2 \cdot \frac{1}{n} \cdot \left(-\frac{1}{K-1}\right) = \frac{2}{n} \cdot \frac{K}{K-1}. \quad (18)$$

217 Let $z = (z_2, \dots, z_K)^T$ and $z^* = (z_2^*, \dots, z_K^*)^T$, with $z_j = \frac{\bar{x}_1 - \bar{x}_j}{\sigma_{1j}}$ and $z_j^* = \frac{\bar{x}_1^* - \bar{x}_j^*}{\sigma_{1j}}$, $j = 2, \dots, K$. We have

$$\begin{aligned} \mathbb{V}(z_j) &= 1, \\ \text{Cor}(z_j, z_k) &= \text{Cov}(\bar{x}_1 - \bar{x}_j, \bar{x}_1 - \bar{x}_k) / (\sigma_{1j}\sigma_{1k}) \\ &= [\mathbb{V}(\bar{x}_1) - 2\text{Cov}(\bar{x}_1, \bar{x}_j) + \text{Cov}(\bar{x}_j, \bar{x}_k)] / (\sigma_{1j}\sigma_{1k}) \\ &= \frac{1}{n} \left(1 + \frac{1}{K-1}\right) / \left(\frac{2}{n} \frac{K}{K-1}\right) = \frac{1}{2}. \end{aligned} \quad (19)$$

218 Thus $z \sim \mathbb{N}(0, \Sigma_0)$ and $z^*|x \sim \mathbb{N}(z, \Sigma_0)$, where Σ_0 is a $(K - 1) \times (K - 1)$ variance matrix with 1 on
 219 the diagonal and $\frac{1}{2}$ on the off-diagonal. Thus

$$P_1 = \mathbb{P}(z_2^* > 0, \dots, z_K^* > 0 | x) = \Phi(z_2, \dots, z_K). \quad (20)$$

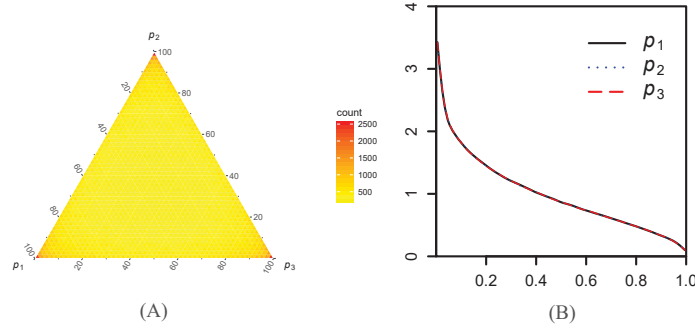


Fig. 2. Marginal and joint distributions of P_1, P_2, P_3 for problem 5 (the multivariate normal example with $K = 3$). The three corners in the plots correspond to points $(1, 0, 0)$, $(0, 1, 0)$, and $(0, 0, 1)$, while the center is $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$. The number of replicates is $R = 10^6$, with $n = 10^6$ and $B = 10^3$.

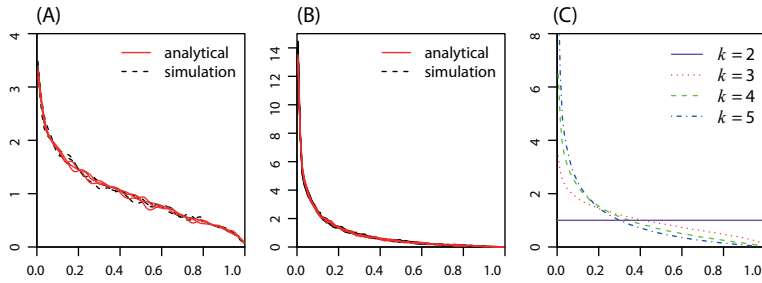


Fig. 3. Marginal distribution of P_1 in comparisons of K equally right or equally wrong and indistinct models based on the normal distribution of Problem 5 ($K = 3$ in A and 6 in B). The sample size is $n = 10^4$. The number of simulated replicates is $R = 10^4$, with $B = 10^3$, but the ‘theoretical’ distribution is based on simulating 10^6 replicates and using Equation 21.

220 As $\bar{x}_j^* - \bar{x}_k^* = (\bar{x}_1^* - \bar{x}_k^*) - (\bar{x}_1^* - \bar{x}_j^*)$, the bootstrap probabilities for all K models given data
 221 x are

$$\begin{bmatrix} P_1 \\ P_2 \\ \vdots \\ P_K \end{bmatrix} = \begin{bmatrix} \Phi(z_2, z_3, \dots, z_K) \\ \Phi(-z_2, z_3 - z_2, \dots, z_K - z_2) \\ \vdots \\ \Phi(-z_K, z_2 - z_K, \dots, z_{K-1} - z_K) \end{bmatrix} \quad (21)$$

222 For example, in the case of $K = 3$, a fast way of simulating the limiting distribution of
 223 (P_1, P_2, P_3) is thus to generate $(z_2, z_3) \sim \mathbb{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}\right)$ and then calculate (P_1, P_2, P_3) by
 224 Equation 21. This is confirmed by the slow simulation of generating x and then x^* in Fig. 2. The
 225 joint distribution of (P_1, P_2, P_3) has peaks at the three corners, and is nearly flat around the center.
 226 By symmetry P_1 has mean $\frac{1}{3}$, and by numerical integration using Equation 20, P_1 has SD =
 227 0.25904. The probability that one of the models is strongly supported is close to 0 (table 1).
 228 Fig. 3A&B shows the marginal distribution of P_1 when $K = 3$ and 6.

229 Even though (P_1, P_2, P_3) do not converge to the point value $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$, extreme bootstrap
 230 support values are not highly frequent. Bootstrap probabilities are thus qualitatively different
 231 from Bayesian model probabilities.

Table 1. Proportions of data replicates with very high bootstrap probability (P_1) in the multivariate normal example (Problem 5)

K	2	3	4	5
$\mathbb{P}\{P_1 > 0.90\}$	0.100	0.023	0.008	0.004
$\mathbb{P}\{P_1 > 0.95\}$	0.050	0.008	0.003	0.001
$\mathbb{P}\{P_1 > 0.99\}$	0.010	0.001	0.000	0.000

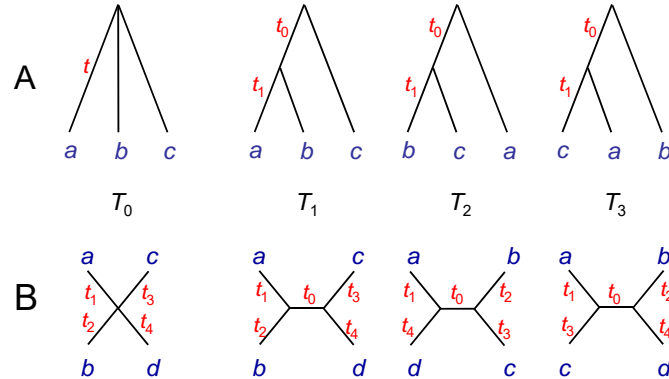


Fig. 4. The star tree T_0 and three binary rooted trees T_1, T_2 , and T_3 for (A) three or (B) four species. Branch length parameters, shown next to the branches, are measured by the expected number of changes per site. The star tree is used to generate data, which are analyzed by ML to compare the three binary trees.

2. BOOTSTRAP IN PHYLOGENETICS

We consider ML reconstruction of phylogenies of three or four species (Fig. 4), under the JC model (Jukes and Cantor, 1969). We simulate data to verify the asymptotic theory and compare with Bayesian results from Yang and Zhu (2018).

Case A (Fig. 5A&A') involves equally right models. This is the star-tree paradox analyzed previously (Lewis *et al.*, 2005; Yang and Rannala, 2005; Yang, 2007a; Susko, 2008). We use the rooted star tree T_0 for three species with $t = 0.2$ (Fig. 4A) to generate datasets to compare the three binary trees. The JC model (Jukes and Cantor, 1969) is used both to generate and to analyze the data. The molecular clock (rate constancy over time) is assumed as well, so that the parameters in each binary tree are the two node ages (t_0 and t_1), measured by the expected number of nucleotide changes per site. The best-fitting parameter values are $t_{0*} = 0$ and $t_{1*} = 0.2$ for each of the three binary trees, so that the three binary trees are equally right models.

Case B (Fig. 5B&B') involves equally wrong models that are indistinct. This is similar to case A except that the JC+ Γ model (Jukes and Cantor, 1969; Yang, 1993) is used to generate data, with different sites in the sequence evolving at variable rates according to the gamma distribution with shape parameter $\alpha = 1$. The data are then analyzed using JC (equivalently to JC+ Γ with $\alpha = \infty$), giving $t_{0*} = 0$ and $t_{1*} = 0.16441$ as the pseudo-true parameter values for each binary tree. The binary trees are equally wrong and indistinct models ($D_1 = D_2 = D_3 > 0$).

Case C (Fig. 5C&C') involves equally wrong and distinct models. Like case B, the simulation model is JC+ Γ with $\alpha = 1$ and the analysis model is JC. However, the molecular clock is not assumed and unrooted trees are used. The true tree is the unrooted star tree T_0 of Fig. 4B, with $t_1 = t_2 = t_3 = t_4 = 0.2$, with $t_{0*} = 0.01037$ and $t_{i*} = 0.16409$, $i = 1, \dots, 4$ for the binary trees (Fig. 4B). As $t_{0*} > 0$, the three binary trees are equally wrong and distinct models ($D_1 = D_2 = D_3 > 0$).

In cases A and B, the data for three species have a multinomial distribution with five

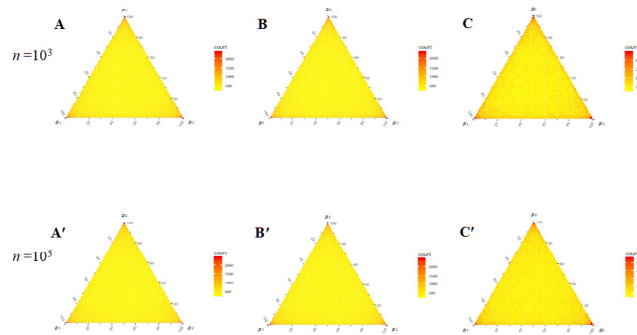


Fig. 5. The joint distribution of the bootstrap model probabilities in the star-tree problem. The star tree T_0 of Fig. 4 is used to simulate data (sequence alignments of $n = 10^3$ or 10^5 sites), and ML is used to compare the three binary trees T_1, T_2 , and T_3 to calculate their bootstrap probabilities (P_1, P_2, P_3) . In (A) and (A'), the true tree is the star tree T_0 for three species of Fig. 4A, with $t = 0.2$. Both the simulation and analysis models are JC, and the three binary trees are equally right models. In (B) and (C), the true tree is the star tree T_0 for three species of Fig. 4A, with $t = 0.2$. The simulation model is JC+ Γ (with $\alpha = 1$), and the analysis model is JC. The three binary trees represent equally wrong and indistinct models. In (C) and (C'), the true tree is the star tree T_0 for four species of Fig. 4B, with $t_1 = t_2 = t_3 = t_4 = 0.2$. The simulation model is JC+ Γ ($\alpha = 1$) and the analysis model is JC. The three binary trees represent equally wrong and distinct models. The number of bootstrap samples $B = 1000$ and the number of replicates is $R = 10^6$ for three-species trees and 10^4 for four-species trees.

257 categories corresponding to the five site patterns xxx , xyx , xyx , yxx , and xyz , where x, y, z are any
 258 distinct nucleotides. Let the frequencies of the informative site patterns xyx , xyx , yxx be \bar{x}_1 , \bar{x}_2 , and
 259 \bar{x}_3 , while that for the two uninformative patterns xxx and xyz be \bar{x}_0 . With the star tree being the
 260 true tree, the probabilities for the three informative site patterns are identical, with $p_1 = p_2 = p_3$.
 261 Tree 1 specifies $p_1 > p_2 = p_3$. Given data x , tree j is the ML tree if \bar{x}_j is the greatest among
 262 \bar{x}_1, \bar{x}_2 , and \bar{x}_3 (Yang, 2000). Then $\bar{x} = (\bar{x}_1, \bar{x}_2, \bar{x}_3)$ is approximately normal, with mean (p, p, p) ,
 263 and variance $p(1-p)/n$ and covariance $-p^2/n$. Applying a multivariate normal approximation
 264 to the multinomial distribution, we see that the problem has the same mathematical structure as
 265 problem 5. Thus the bootstrap distribution for cases A and B should be identical to that in
 266 problem 5. We wrote a C program to simulate and analyze data for cases A and B. Given branch
 267 lengths t_0 and t_1 , the probabilities for the five site patterns are calculated according to the JC
 268 model (Yang, 1994), and the data x are then generated by sampling from the multinomial
 269 distribution. Given data x , bootstrap dataset x^* is sampled using the observed site-pattern
 270 frequencies in x . Then tree j is the ML tree for data x^* if \bar{x}_j^* is the largest among $(\bar{x}_1^*, \bar{x}_2^*, \bar{x}_3^*)$.

271 In case C for four species, the informative site patterns are $xxyy$, $xyxy$ and $xyyx$ while there
 272 are 11 uninformative patterns. The binary tree has only five parameters, such that the model
 273 achieves a better fit to the observed data by having a positive internal branch length. As a result,
 274 the three binary trees are distinct models (with $t_{0*} > 0$). Case C thus differs from problem 5, but
 275 has a similar symmetry in that the K-L distance between any pair of models is the same. From the
 276 general theory, the distribution of bootstrap probabilities (P_1, P_2, P_3) is the same as that in problem
 277 5. We simulated data using EVOLVER, and generated bootstrap resample data using SEQBOOT.
 278 The data are then analyzed using BASEML in PAML (Yang, 2007b).

279 Our theory predicts that the limiting distribution is the same in all three cases, with the
 280 mean $1/3$ and SD 0.25904 . This is confirmed by the simulation (Fig. 5), which gave the mean of
 281 P_1 as $1/3$ and the SD as 0.259 . The bootstrap probabilities have modes at the corners, and
 282 roughly uniformly distributed around the center. While in case C, the Bayesian posterior
 283 probabilities show extreme polarized behavior, concentrated on three points: $(1, 0, 0)$, $(0, 1, 0)$,
 284 and $(0, 0, 1)$ (Yang and Zhu, 2018, Fig. 4C&4C'), bootstrap probabilities are much more
 285 moderate and have a nondegenerate distribution.

Table 2. Proportions of datasets with extreme bootstrap or posterior (in parentheses) probabilities for the three binary trees in the star-tree simulation

n	$\mathbb{P}\{P_{\min} < 1\%\}$	$\mathbb{P}\{P_{\min} < 5\%\}$	$\mathbb{P}\{P_{\max} > 95\%\}$	$\mathbb{P}\{P_{\max} > 99\%\}$	$\mathbb{E}(P_{\min})$	$\mathbb{E}(P_{\max})$
10^3	0.119 (0.234)	0.391 (0.550)	0.028 (0.205)	0.001 (0.079)	0.094 (0.067)	0.644 (0.754)
10^4	0.123 (0.812)	0.386 (0.931)	0.030 (0.606)	0.002 (0.450)	0.093 (0.011)	0.653 (0.897)
10^5	0.113 (0.979)	0.383 (0.992)	0.029 (0.853)	0.004 (0.773)	0.093 (0.001)	0.647 (0.964)

Note.— $P_{\min} = \min(P_1, P_2, P_3)$ and $P_{\max} = \max(P_1, P_2, P_3)$. Data are generated under JC+ Γ with $\alpha = 1$, using the star tree for four species ($a : 0.2, b : 0.2, c : 0.2, d : 0.2$), and analyzed under JC. The number of replicates is $R = 10^3$ and the number of bootstrap samples is $B = 10^3$. The probability density of (P_1, P_2, P_3) is shown in Fig. 5C&C' for $n = 10^3$ and 10^5 , respectively. Posterior tree probabilities from the Bayesian analysis are shown in parentheses, from Yang and Zhu (2018, table S1).

Table 3. Proportions of datasets with strong bootstrap or posterior (in parentheses) support for wrong trees in simulated datasets for four species

n	$\mathbb{P}\{P_1 < 1\%\}$	$\mathbb{P}\{P_1 < 5\%\}$	$\mathbb{P}\{P_{23} > 95\%\}$	$\mathbb{P}\{P_{23} > 99\%\}$
10^3	0.031 (0.083)	0.109 (0.225)	0.019 (0.113)	0.002 (0.038)
10^4	0.009 (0.250)	0.044 (0.337)	0.005 (0.266)	0.000 (0.166)
10^5	0.000 (0.102)	0.001 (0.120)	0.000 (0.115)	0.000 (0.097)

Note.— P_1 is the probability for the true tree, while P_2 and P_3 are for the two wrong trees, with $P_{23} = \max\{P_2, P_3\}$. Data were generated under JC+ Γ with $\alpha = 1$ on the unrooted tree T_1 for four species: $((a : 0.2, b : 0.2) : 0.002, c : 0.2, d : 0.2)$, and analyzed under JC. The number of simulated replicates is $R = 10^3$, with $B = 10^3$. Posterior tree probabilities from the Bayesian analysis are shown in parentheses, from Yang and Zhu (2018, table S2).

286 We also calculated the proportions of datasets in which the bootstrap and posterior
 287 probabilities for the three binary trees are extremely high (table 2). When the sequence length is
 288 $n = 10^5$, $\mathbb{E}(P_{\max}) = 0.647$ using bootstrap method and 0.964 for the Bayesian method. If
 289 $P_{\max} > 0.95$, one of the models is strongly favored, and this occurs in 2.9% of datasets for the
 290 bootstrap and 85.3% for the Bayesian. In other words, it is much less likely to see high bootstrap
 291 support for equally wrong models than high posterior probabilities for them.

292 DISCUSSION

293 As mentioned in Introduction, the interpretation of bootstrap in model selection in general and in
 294 phylogenetics in particular is controversial. A number of studies have attempted to give bootstrap
 295 a Bayesian interpretation, that is, the bootstrap probability for a tree is the probability that the tree
 296 is correct. For example, Hastie *et al.* (2009, p.272) wrote that “[i]n this sense, the bootstrap
 297 distribution represents an (approximate) nonparametric, noninformative posterior distribution for
 298 our parameter.” The plug-in principle for bootstrap appears to support this interpretation:
 299 bootstrap probability $\mathbb{P}\{\Delta^* > 0|x\}$ is an estimate of $\mathbb{P}\{\Delta > 0\}$, which is the probability that the
 300 ML tree is correct. In phylogenetics, such an interpretation was suggested by Efron *et al.* (1996),
 301 although the assumed prior for the corresponding Bayesian analysis has infinite branch lengths
 302 and is implausible biologically (Yang, 2014, p.176).

303 Our analysis suggests qualitatively different asymptotic behaviors between bootstrap and

304 posterior probabilities for models or trees. The greatest difference occurs in the case of
305 comparing equally wrong and distinct models. In that case, the posterior model probabilities
306 show extreme polarized behavior, with $\sim 100\%$ for one model and 0 for others. This behavior is
307 because the log marginal likelihood ratio for two models (or the log Bayes factor) (Δ) is
308 dominated by a random-walk that deviates from 0 (which corresponds to the posterior probability
309 $\frac{1}{2}$ for each model) at the rate of \sqrt{n} when n increases Yang and Zhu (2018), so that for large n , the
310 posterior model probability is either 0 or 1. Bootstrap probabilities show a different behavior.
311 While the log likelihood ratio for the bootstrap dataset (Δ^*) also increases like a random walk
312 when n increases, it is compared with the log likelihood ratio for the original dataset (Δ) when the
313 bootstrap model probability is calculated. As a result, whether the models are distinct or
314 indistinct does not matter anymore.

ACKNOWLEDGMENTS

315
316 We thank Professor Zhi-ming Ma for discussions. This study has been supported by
317 Biotechnology and Biological Sciences Research Council grant (BB/P006493/1) to Z.Y. and a
318 BBSRC equipment grant (BB/R01356X/1). J.H.'s visit to UCL is supported by China
319 Scholarship Council (CSC).

SUPPLEMENTAL DATA

320
321 Data available from the Dryad Digital Repository: <https://doi.org/10.5061/dryad.7m0cfxprw>.

REFERENCES

- 322
323 Berry, V. and Gascuel, O. 1996. On the interpretation of bootstrap trees: appropriate threshold of
324 clade selection and induced gain. *Mol. Biol. Evol.*, 13: 999–1011.
- 325 Bickel, P. J. and Freedman, D. A. 1981. Some asymptotic theory for the bootstrap. *Ann. Statist.*,
326 9: 1196–1217.
- 327 Chan, K. O., Hutter, C. R., Wood, P. L., J., Grismer, L. L., and Brown, R. M. 2020. Larger,
328 unfiltered datasets are more effective at resolving phylogenetic conflict: Introns, exons, and
329 uces resolve ambiguities in golden-backed frogs (anura: Ranidae; genus hylarana). *Mol.*
330 *Phylogenet. Evol.*, 151: 106899.
- 331 Cheng, G. and Huang, J. Z. 2010. Bootstrap consistency for general semiparametric
332 M-estimation. *Ann. Statist.*, 38(5): 2884–2915.
- 333 DasGupta, A. 2008. The bootstrap. In *Asymptotic Theory of Statistics and Probability*, pages
334 461–497. Springer, New York.
- 335 Davison, A. and Hinkley, D. 1997. *Bootstrap Methods and their Application*. Cambridge
336 University Press, Cambridge, UK.
- 337 Dawid, A. 2011. Posterior model probabilities. In P. S. Bandyopadhyay and M. Forster, editors,
338 *Philosophy of Statistics*, pages 607–630. Elsevier, New York.
- 339 Efron, B. 1979. Bootstrap methods: another look at the jackknife. *Ann. Stat.*, 7: 1–26.
- 340 Efron, B. and Tibshirani, R. 1993. *An Introduction to the Bootstrap*. Chapman and Hall, London.
- 341 Efron, B., Halloran, E., and Holmes, S. 1996. Bootstrap confidence levels for phylogenetic trees.
342 *Proc. Natl. Acad. Sci. U.S.A.*, 93: 13429–13434.

- 343 Felsenstein, J. 1981. Evolutionary trees from DNA sequences: a maximum likelihood approach.
344 *J. Mol. Evol.*, 17(6): 368–376.
- 345 Felsenstein, J. 1985. Confidence limits on phylogenies: an approach using the bootstrap.
346 *Evolution*, 39: 783–791.
- 347 Felsenstein, J. and Kishino, H. 1993. Is there something wrong with the bootstrap on
348 phylogenies? a reply to Hillis and Bull. *Syst. Biol.*, 42(2): 193–200.
- 349 Fitch, W. M. 1971. Toward defining the course of evolution: minimum change for a specific tree
350 topology. *Syst. Zool.*, 20: 406–416.
- 351 Hastie, T., Tibshirani, R., and Friedman, J. 2009. *The Elements of Statistical Learning: Data*
352 *Mining, Inference, and Prediction*. Springer, New York, 2 edition.
- 353 Hillis, D. M. and Bull, J. J. 1993. An empirical test of bootstrapping as a method for assessing
354 confidence in phylogenetic analysis. *Syst. Biol.*, 42: 182–192.
- 355 Holmes, S. 2003. Bootstrapping phylogenetic trees: theory and methods. *Stat. Sci.*, 18: 241–255.
- 356 Huelsenbeck, J. and Rannala, B. 2004. Frequentist properties of bayesian posterior probabilities
357 of phylogenetic trees under simple and complex substitution models. *Syst. Biol.*, 53: 904–913.
- 358 Huggins, J. H. and Miller, J. W. 2020. Using bagged posteriors for robust inference and model
359 criticism. page arXiv:1912.07104.
- 360 Jukes, T. and Cantor, C. 1969. Evolution of protein molecules. In H. Munro, editor, *Mammalian*
361 *Protein Metabolism*, pages 21–123. Academic Press, New York.
- 362 Lemoine, F., Domelevo Entfellner, J.-B., Wilkinson, E., Correia, D., Davila Felipe, M.,
363 De Oliveira, T., and Gascuel, O. 2018. Renewing Felsenstein’s phylogenetic bootstrap in the
364 era of big data. *Nature*, 556: 452–456.
- 365 Lewis, P., Holder, M., and Holsinger, K. 2005. Polytomies and bayesian phylogenetic inference.
366 *Syst. Biol.*, 54: 241–253.
- 367 O’Hagan, A. and Forster, J. 2004. *Kendall’s Advanced Theory of Statistics: Bayesian Inference*.
368 Arnold, London.
- 369 Rannala, B. and Yang, Z. 1996. Probability distribution of molecular evolutionary trees: a new
370 method of phylogenetic inference. *J. Mol. Evol.*, 43: 304–311.
- 371 Rubin, D. B. 1981. The Bayesian bootstrap. *Ann. Statist.*, 9(1): 130–134.
- 372 Saitou, N. and Nei, M. 1987. The neighbor-joining method: a new method for reconstructing
373 phylogenetic trees. *Mol. Biol. Evol.*, 4: 406–425.
- 374 Susko, E. 2008. On the distributions of bootstrap support and posterior distributions for a star
375 tree. *Syst. Biol.*, 57: 602–612.
- 376 Susko, E. 2009. Bootstrap support is not first-order correct. *Syst. Biol.*, 58(2): 211–223.
- 377 Susko, E. 2010. First-order correct bootstrap support adjustments for splits that allow hypothesis
378 testing when using maximum likelihood estimation. *Mol. Biol. Evol.*, 27: 1621–1629.
- 379 Weng, C.-S. 1989. On a second-order asymptotic property of the bayesian bootstrap mean. *Ann.*
380 *Statist.*, 17(2): 705–710.

- 381 White, H. 1982. Maximum likelihood estimation of misspecified models. *Econometrica*, 50:
382 1–25.
- 383 Yang, Z. 1993. Maximum-likelihood estimation of phylogeny from DNA sequences when
384 substitution rates differ over sites. *Mol. Biol. Evol.*, 10: 1396–1401.
- 385 Yang, Z. 1994. Statistical properties of the maximum likelihood method of phylogenetic
386 estimation and comparison with distance matrix methods. *Syst. Biol.*, 43: 329–342.
- 387 Yang, Z. 2000. Complexity of the simplest phylogenetic estimation problem. *Proc. R. Soc. B:*
388 *Biol. Sci.*, 267: 109–116.
- 389 Yang, Z. 2007a. Fair-balance paradox, star-tree paradox and bayesian phylogenetics. *Mol. Biol.*
390 *Evol.*, 24: 1639–1655.
- 391 Yang, Z. 2007b. PAML 4: Phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.*, 24:
392 1586–1591.
- 393 Yang, Z. 2014. *Molecular Evolution: A Statistical Approach*. Oxford University Press, Oxford,
394 England.
- 395 Yang, Z. and Rannala, B. 2005. Branch-length prior influences bayesian posterior probability of
396 phylogeny. *Syst. Biol.*, 54: 455–470.
- 397 Yang, Z. and Zhu, T. 2018. Bayesian selection of misspecified models is overconfident and may
398 cause spurious posterior probabilities for phylogenetic trees. *Proc. Natl. Acad. Sci. USA*,
399 115(8): 1854–1859.
- 400 Zharkikh, A. and Li, W.-H. 1992. Statistical properties of bootstrap estimation of phylogenetic
401 variability from nucleotide sequences. i. four taxa with a molecular clock. *Mol. Biol. Evol.*, 9:
402 1119–1147.
- 403 Zharkikh, A. and Li, W.-H. 1995. Estimation of confidence in phylogeny: the
404 complete-and-partial bootstrap technique. *Mol. Phylogenet. Evol.*, 4: 44–63.

405 APPENDIX. ASYMPTOTIC THEORY FOR BOOTSTRAP PROBABILITY IN MODEL SELECTION

406 We use ML to compare K models, $H_j : X \sim f_j(X|\theta_j)$, $j = 1, \dots, K$. The dataset, $x = \{x_1, \dots, x_n\}$,
 407 is an i.i.d. sample of from the true model $g(X)$. Given x , we generate a bootstrap sample x^* and
 408 analyze it using ML. The bootstrap probability for model H_1 is the probability that model H_1 has
 409 higher log likelihood than other models in the bootstrap sample.

410 **The case of two equally wrong and distinct models.** We decompose the log-likelihood
 411 ratio between models H_1 and H_2 for the bootstrap dataset x^* into several components, and study
 412 their dynamics when $n \rightarrow \infty$.

$$\Delta^* \equiv \log \frac{f_1(x^*|\hat{\theta}_1^*)}{f_2(x^*|\hat{\theta}_2^*)} = \log \frac{f_1(x^*|\hat{\theta}_1^*)}{f_1(x^*|\hat{\theta}_1)} - \log \frac{f_2(x^*|\hat{\theta}_2^*)}{f_2(x^*|\hat{\theta}_2)} + \log \frac{f_1(x^*|\hat{\theta}_1)}{f_2(x^*|\hat{\theta}_2)} \equiv \Delta A_1 - \Delta A_2 + \Delta_*^*. \quad (\text{A1})$$

413 Model H_1 is the selected model in the bootstrap sample if and only if $\Delta^* > 0$, so that the bootstrap
 414 probability for H_1 given data x is $P_1 \equiv \mathbb{P}\{\Delta^* > 0|x\}$. We are interested in the distribution of P_1
 415 when x varies. First we consider the case where H_1 and H_2 are equally wrong and distinct. We
 416 show that ΔA_1 and ΔA_2 are $O_p(1)$ while Δ_*^* is $O_p(n^{1/2})$, so that Δ^* is dominated by Δ_*^* .

417 Taking the same approach as in Dawid (2011) and Yang and Zhu (2018), we apply Taylor
 418 expansion to the log likelihood, $\log f_1(x^*|\theta_1)$, for the bootstrap dataset x^* around the MLE $\hat{\theta}_1^*$ and
 419 then let $\theta_1 = \hat{\theta}_1$. We have

$$\begin{aligned} \Delta A_1 &= \log f_1(x^*|\hat{\theta}_1^*) - \log f_1(x^*|\hat{\theta}_1) \\ &\approx \frac{1}{2} \{(\hat{\theta}_1^* - \hat{\theta}_1)\}^T \{nJ_1(\hat{\theta}_1^*)\} \{(\hat{\theta}_1^* - \hat{\theta}_1)\} \\ &\approx \frac{1}{2} \{\sqrt{n}(\hat{\theta}_1 - \theta_{1*})\}^T J_1(\theta_{1*}) \{\sqrt{n}(\hat{\theta}_1 - \theta_{1*})\}, \end{aligned} \quad (\text{A2})$$

420 where $J_1(\theta_1) = \mathbb{E}\{-\nabla^2 \log f_1(X|\theta_1)\}$ and ∇^2 is the second derivatives with respect to θ_1 . From
 421 the plug-in principle, x^* varies given $\hat{\theta}$ as does x given θ_* (Efron and Tibshirani, 1993). We have
 422 $\sqrt{n}(\hat{\theta}_1^* - \hat{\theta}_1) \xrightarrow{d} \sqrt{n}(\hat{\theta}_1 - \theta_{1*})$ (Bickel and Freedman, 1981; Cheng and Huang 2010, Theorem
 423 2), and

$$\sqrt{n}(\hat{\theta}_1 - \theta_{1*}) \sim \mathbb{N}(0, [J_1(\theta_{1*})^{-1}]^T I_1(\theta_{1*}) J_1(\theta_{1*})^{-1}), \quad (\text{A3})$$

424 where $I_1(\theta_1) = \mathbb{E}\{\nabla \log f_1(X|\theta_1) \cdot \nabla \log f_1(X|\theta_1)^T\}$ (White, 1982, Theorem 3.2). Thus ΔA_1 is a
 425 quadratic form of normal variates and is $O_p(1)$. If H_1 is the true model, $\Delta A_1 \sim \frac{1}{2} \chi_d^2$ where d is the
 426 number of parameters in H_1 . Similarly $\Delta A_2 = O_p(1)$.

427 We write the third term in Equation A1 as

$$\Delta_*^* \equiv \log \frac{f_1(x^*|\hat{\theta}_1)}{f_2(x^*|\hat{\theta}_2)} = \sum_{i=1}^n \log \frac{f_1(x_i^*|\hat{\theta}_1)}{f_2(x_i^*|\hat{\theta}_2)} \equiv \sum_{i=1}^n r_i^*(x). \quad (\text{A4})$$

428 Define two log-likelihood ratios based on the original data x ,

$$\begin{aligned} \Delta_* &\equiv \log \frac{f_1(x|\theta_{1*})}{f_2(x|\theta_{2*})}, \\ \Delta &\equiv \log \frac{f_1(x|\hat{\theta}_1)}{f_2(x|\hat{\theta}_2)} = \sum_{i=1}^n \log \frac{f_1(x_i|\hat{\theta}_1)}{f_2(x_i|\hat{\theta}_2)} \equiv \sum_{i=1}^n r_i. \end{aligned} \quad (\text{A5})$$

429 Note that Δ_* is a sum of n i.i.d. terms, so that when $n \rightarrow \infty$, $\Delta_* \sim \mathbb{N}(0, n\sigma^2)$, with $\mathbb{E}(\Delta_*) =$

430 $n(D_1 - D_2) = 0$ (Equation 1) and $\mathbb{V}(\Delta_*) = n\sigma^2$, where

$$\sigma^2 \equiv \mathbb{V}_g \left\{ \log \frac{f_1(X|\theta_{1*})}{f_2(X|\theta_{2*})} \right\} = \int g(X) \left[\log \frac{f_1(X|\theta_{1*})}{f_2(X|\theta_{2*})} \right]^2 dX. \quad (\text{A6})$$

431 When $n \rightarrow \infty$, $\bar{r} = \frac{1}{n} \sum_{i=1}^n r_i \rightarrow D_1 - D_2 = 0$ and $s^2 = \frac{1}{n} \sum_{i=1}^n (r_i - \bar{r})^2 \rightarrow \sigma^2$, so that $\Delta \sim \mathbb{N}(0, n\sigma^2)$.
432 Given data x , $\{r_i^*\}$ are conditionally independent, with expectation and variance

$$\begin{aligned} \mathbb{E}(\Delta_*^*|x) &= n \mathbb{E} \left\{ \log \frac{f_1(x_1^*|\hat{\theta}_1)}{f_2(x_1^*|\hat{\theta}_2)} \middle| x \right\} \approx n \cdot \frac{1}{n} \sum_{i=1}^n \log \frac{f_1(x_i|\hat{\theta}_1)}{f_2(x_i|\hat{\theta}_2)} = \sum_{i=1}^n r_i = \Delta, \\ \mathbb{V}(\Delta_*^*|x) &= n \mathbb{V} \left\{ \log \frac{f_1(x_1^*|\hat{\theta}_1)}{f_2(x_1^*|\hat{\theta}_2)} \middle| x \right\} = n \mathbb{E} \{ (r_i^* - \mathbb{E}(r_i^*))^2 | x \} \approx n\sigma^2. \end{aligned} \quad (\text{A7})$$

433 Thus $\Delta_*^*|x \sim \mathbb{N}(\Delta, n\sigma^2)$. The bootstrap probability for H_1 is

$$P_1 = \mathbb{P}\{\Delta^* > 0|x\} = \mathbb{P}\{\Delta A_1 - \Delta A_2 + \Delta_*^* > 0|x\} \approx \mathbb{P}\{\Delta_*^* > 0|x\} \approx \Phi\left(\frac{\Delta}{\sqrt{n\sigma}}\right) \sim \mathbb{U}(0, 1). \quad (\text{A8})$$

434 P_1 varies among datasets like a random number.

435 The case where there are no free parameters in the compared models has been discussed
436 in the main paper. We have

$$\Delta = \Delta_* = \log \frac{f_1(x)}{f_2(x)}, \quad \Delta^* = \Delta_*^* = \log \frac{f_1(x^*)}{f_2(x^*)}, \quad (\text{A9})$$

437 with $\Delta \sim \mathbb{N}(0, n\sigma^2)$ and $\Delta^* \sim \mathbb{N}(\Delta, n\sigma^2)$, as $n \rightarrow \infty$. Thus

$$P_1 = \mathbb{P}\{\Delta^* > 0|x\} = \Phi\left(\frac{\Delta}{\sqrt{n\sigma}}\right) \rightarrow \mathbb{U}(0, 1). \quad (\text{A10})$$

438 The case where the two models are equally right or are equally wrong and indistinct, that
439 is, with $f_1(X|\theta_1^*) = f_2(X|\theta_2^*)$ for almost every X . We have $\Delta_* = 0$ in Equation A5, and
440 $\Delta = O_p(1)$. As a result, $\Delta_*^* = O_p(1)$, as well as $\Delta A_1 = O_p(1)$ and $\Delta A_2 = O_p(1)$. From Equation
441 A2, ΔA_1 and ΔA_2 have the same distribution, with $\mathbb{E}(\Delta A_1 - \Delta A_2|x) = 0$. Thus $\mathbb{E}(\Delta^*|x) = \mathbb{E}(\Delta_*^*|x)$
442 $= \Delta$. Let F be the CDF of Δ , which has mean 0. Then

$$P_1 = \mathbb{P}\{\Delta^* > 0|x\} = 1 - F(-\Delta). \quad (\text{A11})$$

443 Thus with $n \rightarrow \infty$, P_1 converges to a non-degenerate distribution, which is $\mathbb{U}(0, 1)$ if and only if
444 $\Delta^* - \Delta$ has the same distribution as $-\Delta$.

445 DasGupta (2008, Chapter 29) discusses regularity conditions under which $T^* - T$ and
446 $T - \mathbb{E}(T)$ have the same distribution, so that the bootstrap plugin principle can be applied, where
447 T is a statistic or function of data x . If those conditions are not satisfied, the standard bootstrap
448 will fail as $T^* - T$ will not approximate $T - \mathbb{E}(T)$. Problem 4 is one such case, and $\Delta^* - \Delta$ and Δ
449 have different distributions, and the limiting distribution of P_1 is not uniform. As indistinct
450 models are more similar to each other than distinct models and as $P_1 \sim \mathbb{U}(0, 1)$ when the two
451 models are distinct (and equally wrong), we conjecture that $\mathbb{V}(P_1) \leq \frac{1}{12}$, the variance of $\mathbb{U}(0, 1)$.

452 Problems 3 and 4 are examples of equally right or equally wrong but indistinct models.
453 Problem 3 shows the $\mathbb{U}(0, 1)$ distribution, while problem 4 shows a non-uniform distribution.

454 **The case of K models.** Let the K models be H_1, \dots, H_K , all of which have the same K-L
455 distance to the true model. Define

$$\Delta_{*jk} \equiv \log \frac{f_j(x|\theta_{*j})}{f_k(x|\theta_{*k})}, \quad \Delta_{jk} \equiv \sum_{i=1}^n \log \frac{f_j(x_i|\hat{\theta}_j)}{f_k(x_i|\hat{\theta}_k)} \quad (\text{A12})$$

456 for dataset x and

$$\Delta_{*jk}^* \equiv \sum_{i=1}^n \log \frac{f_j(x_i^*|\hat{\theta}_j)}{f_k(x_i^*|\hat{\theta}_k)}, \quad \Delta_{jk}^* \equiv \sum_{i=1}^n \log \frac{f_j(x_i^*|\hat{\theta}_j^*)}{f_k(x_i^*|\hat{\theta}_k^*)} \quad (\text{A13})$$

457 for bootstrap dataset x^* .

458 First consider the case where the K models are equally wrong and distinct. As in the case
459 of two models, Δ_{jk}^* is dominated by Δ_{*jk}^* so that $\Delta_{jk}^* \approx \Delta_{*jk}^*$ while $\Delta \sim \mathbb{N}(0, n\sigma_{jk}^2)$ and $\Delta_{*jk}^* \sim$
460 $\mathbb{N}(\Delta_{jk}, n\sigma_{jk}^2)$, with $\sigma_{jk}^2 \equiv \mathbb{V} \left\{ \log \frac{f_j(X|\theta_{j*})}{f_k(X|\theta_{k*})} \right\}$ (see Equation A6). Given x , there will be a set of
461 bootstrap probabilities (P_1, \dots, P_K) . For example

$$P_1 = \mathbb{P}\{\Delta_{12}^* > 0, \dots, \Delta_{1K}^* > 0|x\} \approx \mathbb{P}\{\Delta_{*12}^* > 0, \dots, \Delta_{*1K}^* > 0|x\}. \quad (\text{A14})$$

462 Let $z = \{z_2, \dots, z_{K-1}\}$ and $z^* = \{z_2^*, \dots, z_{K-1}^*\}$, where $z_j = \frac{\Delta_{1j}}{\sqrt{n}\sigma_{1j}}$ and $z_j^* = \frac{\Delta_{1j}^*}{\sqrt{n}\sigma_{1j}}$. Let

$$\rho_{jk} = \text{Cor}(z_j, z_k) = \text{Cor}(\Delta_{1j}, \Delta_{1k}) = \frac{1}{\sigma_{1j}\sigma_{1k}} \text{Cov} \left(\log \frac{f_1(X|\theta_{1*})}{f_j(X|\theta_{j*})}, \log \frac{f_1(X|\theta_{1*})}{f_k(X|\theta_{k*})} \right). \quad (\text{A15})$$

463 Thus $z \sim \mathbb{N}(0, \Sigma_0)$ and $z^*|x \sim \mathbb{N}(z, \Sigma_0)$, where Σ_0 is a $(K-1) \times (K-1)$ variance matrix with 1 on
464 the diagonal and ρ_{jk} on the off-diagonal. We have

$$P_1 = \mathbb{P}(z_2^* > 0, \dots, z_K^* > 0|x) = \Phi(-z_2, \dots, -z_K), \quad (\text{A16})$$

465 where Φ is the $(K-1)$ -variate CDF of $\mathbb{N}(0, \Sigma_0)$. Bootstrap probabilities for the other models,
466 P_2, \dots, P_K , are given similarly.

467 When there is strong symmetry in the problem so that the K-L distance between any two
468 models is the same, the variance matrix Σ_0 will have 1 on the diagonal and $\rho_{jk} = \frac{1}{2}$ on the
469 off-diagonal, and further simplifications are possible. The joint distribution of bootstrap model
470 probabilities (P_1, \dots, P_K) can be simulated as follows (see Problem 5). Sample $z = \{z_2, \dots, z_K\} \sim$
471 $\mathbb{N}(0, \Sigma_0)$ where Σ_0 is $(K-1) \times (K-1)$, with 1 on the diagonal and $\frac{1}{2}$ on the off-diagonal. Let
472 $z_1 = -(z_2 + \dots + z_K)$. Then calculate

$$\begin{bmatrix} P_1 \\ P_2 \\ \vdots \\ P_K \end{bmatrix} = \begin{bmatrix} \Phi(z_2, z_3, \dots, z_K) \\ \Phi(-z_2, z_3 - z_2, \dots, z_K - z_2) \\ \vdots \\ \Phi(-z_K, z_2 - z_K, \dots, z_{K-1} - z_K) \end{bmatrix}. \quad (\text{A17})$$

473 If the K models under comparison are equally right or equally wrong and indistinct, Δ_{jk}^*
474 $= O_p(1)$. Then the bootstrap probabilities (P_1, \dots, P_K) have a nondegenerate distribution.

475 In the case where some of the K models are equally wrong and distinct while others are
476 indistinct, the dynamics of bootstrap support values may be complex. See table S1 for examples.

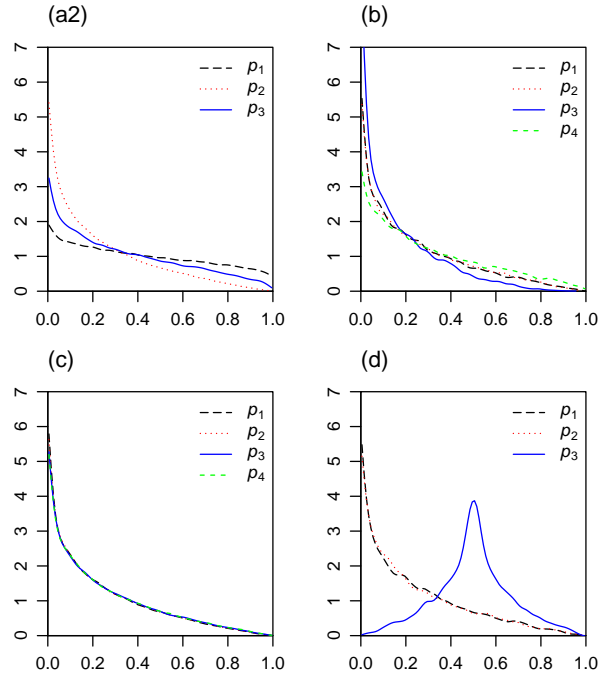


Fig. S1. The density of marginal distribution of bootstrap support P_i , $i = 1, 2, 3, 4$ for cases (a2),(b),(c) and (d).

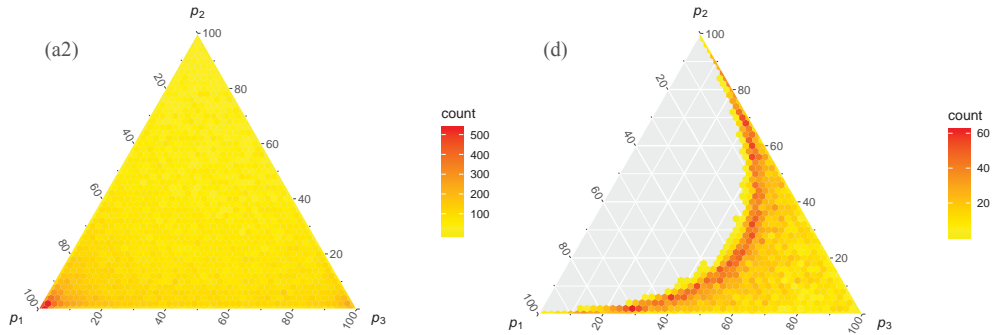


Fig. S2. Ternary distribution of bootstrap support for cases (a2) and (d).

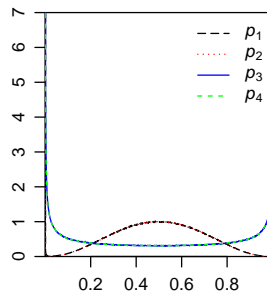


Fig. S3. The density of marginal distribution of posterior probability P_i , $i = 1, 2, 3, 4$ for case (c).

1 SI TEXT. ANALYSIS OF PROBLEM 4: MODELS OF NORMAL MEAN AND VARIANCE

2 The true model is $\mathbb{N}(0, 1)$ and the two compared models are $\mathbb{N}(\mu, 1)$ and $\mathbb{N}(0, \nu)$. In H_1 , $\mu_* = 0$
 3 while in H_2 , $\nu_* = 1$, so the two models are equally right. The data x may be summarized as the
 4 sample mean \bar{x} and sample variance $s^2 = \frac{1}{n} \sum_i (x_i - \bar{x})^2$. The MLE of the parameter is $\hat{\mu} = \bar{x} =$
 5 $\frac{1}{n} \sum x_i$ under H_1 and $\hat{\nu} = \bar{x}^2 = \frac{1}{n} \sum x_i^2 = \bar{x}^2 + s^2$ under H_2 . The log-likelihood values are

$$\begin{aligned} \ell_1(\hat{\mu}) &= -\frac{1}{2} \sum (x_i - \bar{x})^2 = -\frac{1}{2} n s^2, \\ \ell_2(\hat{\nu}) &= -\frac{n}{2} \log \left(\frac{1}{n} \sum x_i^2 \right) - \frac{n}{2} = -\frac{n}{2} \log(s^2 + \bar{x}^2) - \frac{n}{2} \end{aligned} \quad (\text{S1})$$

6 First we derive Equation 15 from first principles. Then we analyze the problem following
 7 the general theory of the Appendix. The problem has been simplified to the following. Let $z_1 =$
 8 $\sqrt{n}\bar{x} \sim \mathbb{N}(0, 1)$ and $z_2 = \sqrt{\frac{n}{2}}(s^2 - 1) \sim \mathbb{N}(0, 1)$, and define z_1^* and z_2^* accordingly, with $z_1^* | z_1 \sim$
 9 $\mathbb{N}(z_1, 1)$ and $z_2^* | z_2 \sim \mathbb{N}(z_2, 1)$. We seek the distribution of P_1 when x varies, defined as

$$P_1 = \mathbb{P}\{|z_1^*| > |z_2^*| \mid z_1, z_2\}. \quad (\text{S2})$$

10 The CDF of the difference between two independent folded normal random variables with
 11 the same variance, $X_1 \sim \mathbb{N}(\mu_1, \sigma^2)$ and $X_2 \sim \mathbb{N}(\mu_2, \sigma^2)$, is

$$\begin{aligned} \mathbb{P}(|X_1| - |X_2| < t) &= \Phi\left(\frac{t}{\sqrt{2}\sigma} - \frac{\tilde{\mu}_1}{\sigma}\right) \left(1 - \Phi\left(\frac{-t}{\sqrt{2}\sigma} - \frac{\tilde{\mu}_2}{\sigma}\right)\right) \\ &+ \left(1 - \Phi\left(\frac{-t}{\sqrt{2}\sigma} - \frac{\tilde{\mu}_1}{\sigma}\right)\right) \Phi\left(\frac{t}{\sqrt{2}\sigma} - \frac{\tilde{\mu}_2}{\sigma}\right), \end{aligned} \quad (\text{S3})$$

12 where $\begin{bmatrix} \tilde{\mu}_1 \\ \tilde{\mu}_2 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} \mu_1 - \mu_2 \\ \mu_1 + \mu_2 \end{bmatrix}$. Thus

$$\begin{aligned} P_1 &= \Phi\left(-\frac{1}{\sqrt{2}}(z_2 - z_1)\right) \left(1 - \Phi\left(-\frac{1}{\sqrt{2}}(z_1 + z_2)\right)\right) \\ &+ \left(1 - \Phi\left(-\frac{1}{\sqrt{2}}(z_2 - z_1)\right)\right) \Phi\left(-\frac{1}{\sqrt{2}}(z_1 + z_2)\right) \\ &\equiv Y_1 Y_2 + (1 - Y_1)(1 - Y_2), \end{aligned} \quad (\text{S4})$$

13 where Y_1 and Y_2 are i.i.d. from $\mathbb{U}(0, 1)$. While P_1 is a function of data x (Equation 14) or of z_1 and
 14 z_2 (Equation S2), it is now considered a function of Y_1 and Y_2 . We use a further variable
 15 transform. Let $W_1 = 2Y_1 - 1$ and $W_2 = 2Y_2 - 1$, with W_1 and W_2 i.i.d. from $\mathbb{U}(-1, 1)$. The density
 16 of $T = W_1 W_2$ is given by the product convolution of W_1 and W_2 as

$$f_T(t) = \int_{\mathbb{R}} f_{W_1}(w_1) f_{W_2}(t/w_1) \frac{1}{|w_1|} dw_1 = \int_{|w_1| < 1, |t/w_1| < 1} \frac{1}{4|w_1|} dw_1 = -\frac{1}{2} \log |t|, \quad (\text{S5})$$

17 $-1 < t < 1$. Note that the region of integral consists of four disjoint intervals: $0 < t < w_1 < 1$,
 18 $0 < t < -w_1 < 1$, $0 < -t < w_1 < 1$, and $0 < -t < -w_1 < 1$. As $P_1 = \frac{1}{2}(T + 1)$, Equation 15
 19 follows.

20 Next we analyze the problem by working with the log likelihood ratios as in the general
 21 theory. The log-likelihoods under the two models are

$$\begin{aligned} \ell_1(\mu) &= -\frac{1}{2} \sum (x_i - \mu)^2 = -\frac{n}{2} (\bar{x}^2 - 2\bar{x}\mu + \mu^2), \\ \ell_2(\nu) &= -\frac{n}{2} \log \nu - \frac{1}{2\nu} \sum x_i^2 = -\frac{n}{2} (\log \nu + \frac{1}{\nu} \bar{x}^2). \end{aligned} \quad (\text{S6})$$

22 At the MLE $\hat{\mu} = \bar{x}$ in H_1 and $\hat{v} = \bar{x}^2$ in H_2 , we have $\ell_1(\hat{\mu}) = -\frac{1}{2}n(\bar{x}^2 - \bar{x}^2)$ and $\ell_2(\hat{v}) =$
 23 $-\frac{n}{2}(\log \bar{x}^2 + 1)$, so that

$$\begin{aligned}\Delta &= \ell_1(\hat{\mu}) - \ell_2(\hat{v}) = \frac{n}{2} [\log \bar{x}^2 + 1 - \bar{x}^2 + \bar{x}^2] \\ &= \frac{n}{2} g(\hat{\mu}, \hat{v}) = \frac{n}{2} [g(\hat{\mu}, \hat{v}) - g(\mu_*, v_*)]\end{aligned}\quad (\text{S7})$$

24 Here we define a function $g(\mu, v) = \log(v) + 1 - v + \mu^2$. We have $\nabla g(\mu, v) = (2\mu, \frac{1}{v} - 1)$ and
 25 $\nabla^2 g(\mu, v) = \begin{pmatrix} 2 & 0 \\ 0 & -1/v^2 \end{pmatrix}$. Note that $g(\mu_*, v_*) = 0$, and $\nabla g(\mu_*, v_*) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ and $\nabla^2 g(\mu_*, v_*) = \begin{pmatrix} 2 & 0 \\ 0 & -1 \end{pmatrix}$.
 26 Let $\bar{Y} = (\bar{x}, \bar{x}^2)$ and $\mu_Y = (\mu_*, v_*)$, with $g(\mu_Y) = \Delta_* = 0$. Then applying Taylor expansion, we get

$$\begin{aligned}\Delta &= \frac{n}{2} [g(\bar{Y}) - g(\mu_Y)] \approx \frac{n}{2} \cdot \frac{1}{2} (\hat{\mu}, \hat{v} - 1) \begin{pmatrix} 2 & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} \hat{\mu} \\ \hat{v} - 1 \end{pmatrix} \\ &= \frac{1}{2} n \bar{x}^2 - \frac{1}{2} \frac{n}{2} (s^2 + \bar{x}^2 - 1)^2 = \frac{1}{2} z_1^2 - \frac{1}{2} z_2^2,\end{aligned}\quad (\text{S8})$$

27 where $z_1 = \sqrt{n}\bar{x} \sim \mathbb{N}(0, 1)$ and $z_2 = \sqrt{\frac{n}{2}}(s^2 + \bar{x}^2 - 1) \approx \sqrt{\frac{n}{2}}(s^2 - 1) \sim \mathbb{N}(0, 1)$, as defined in the
 28 main text.

29 For the bootstrap dataset x^* , we have $\hat{\mu}^* = \bar{x}^*$ and $\hat{v}^* = \bar{x}^{2*} = s^{2*} + \bar{x}^{2*}$. The three terms in
 30 Equation A1 can be derived as follows.

$$\begin{aligned}\Delta A_1 &\equiv \log \frac{f_1(x^* | \hat{\theta}_1^*)}{f_1(x^* | \hat{\theta}_1)} = -\frac{1}{2} \sum (x_i^* - \hat{\mu}^*)^2 + \frac{1}{2} \sum (x_i^* - \hat{\mu})^2 \\ &= -\frac{n}{2} (\bar{x}^{2*} - \bar{x}^{2*} - \bar{x}^{2*} + 2\bar{x}^* \bar{x} - \bar{x}^2) \\ &= \frac{n}{2} (\bar{x}^{2*} - 2\bar{x}^* \bar{x} + \bar{x}^2) = \frac{n}{2} (\bar{x}^* - \bar{x})^2 = \frac{1}{2} (z_1^* - z_1)^2,\end{aligned}\quad (\text{S9})$$

31 with z_1^* and z_2^* defined above. Similarly

$$\begin{aligned}\Delta A_2 &\equiv \log \frac{f_2(x^* | \hat{\theta}_2^*)}{f_2(x^* | \hat{\theta}_2)} = -\frac{n}{2} (\log \hat{v}^* + \frac{\bar{x}^{2*}}{\hat{v}^*}) + \frac{n}{2} (\log \hat{v} + \frac{\bar{x}^{2*}}{\hat{v}}) \\ &= -\frac{n}{2} [(\log \hat{v}^* - \frac{\hat{v}^*}{\hat{v}}) - (\log \hat{v} - 1)] \\ &\approx -\frac{n}{2} \cdot 0 \cdot (\hat{v}^* - \hat{v}) - \frac{n}{2} \cdot \frac{1}{2} (-\frac{1}{\hat{v}^2}) (\hat{v}^* - \hat{v})^2 \approx \frac{n}{4} (\hat{v}^* - \hat{v})^2 = \frac{1}{2} (z_2^* - z_2)^2.\end{aligned}\quad (\text{S10})$$

32 Here in the Taylor expansion \hat{v} is constant and \hat{v}^* is a random variable given x .

$$\begin{aligned}\Delta_*^* &\equiv \log \frac{f_1(x^* | \hat{\theta}_1)}{f_2(x^* | \hat{\theta}_2)} = -\frac{1}{2} \sum (x_i^* - \hat{\mu})^2 + \frac{n}{2} (\log \hat{v} + \frac{\bar{x}^{2*}}{\hat{v}}) = \frac{n}{2} (-\hat{v}^* + 2\bar{x}^* \bar{x} - \bar{x}^2 + \log \hat{v} + \frac{\hat{v}^*}{\hat{v}}) \\ &\approx \frac{n}{2} (2\bar{x}^* \bar{x} - 2\bar{x}^2) + \frac{n}{2} \bar{x}^2 + \frac{n}{2} [(\frac{\hat{v}^*}{\hat{v}} - \hat{v}^*) - (1 - \hat{v})] + \frac{n}{2} (\log \hat{v} + 1 - \hat{v}) \\ &\approx n\bar{x}(\bar{x}^* - \bar{x}) + \frac{n}{2} (\frac{1}{\hat{v}} - 1) (\hat{v}^* - \hat{v}) + \frac{n}{2} (\bar{x}^2 + \log \hat{v} + 1 - \hat{v}) \\ &= n\bar{x}(\bar{x}^* - \bar{x}) - \frac{n}{2} \frac{1}{\hat{v}} (\hat{v} - 1) (\hat{v}^* - \hat{v}) + \Delta \approx z_1 (z_1^* - z_1) - z_2 (z_2^* - z_2) + \Delta.\end{aligned}\quad (\text{S11})$$

33 Thus

$$\begin{aligned}\Delta^* &= \Delta A_1 - \Delta A_2 + \Delta_*^* \\ &\approx \frac{1}{2} (z_1^* - z_1)^2 - \frac{1}{2} (z_2^* - z_2)^2 + z_1 (z_1^* - z_1) - z_2 (z_2^* - z_2) + \frac{1}{2} (z_1^2 - z_2^2) \\ &= \frac{1}{2} (z_1^{*2} - z_2^{*2}).\end{aligned}\quad (\text{S12})$$

34 This can also be obtained by a Taylor expansion of the function $g(\cdot)$,

$$\begin{aligned}\Delta^* - \Delta &= \frac{n}{2}[g(\bar{Y}^*) - g(\bar{Y})] = \frac{n}{2}[g(\hat{\mu}^*, \hat{\nu}^*) - g(\hat{\mu}, \hat{\nu})] \\ &\approx \frac{n}{2}(2\hat{\mu}, \frac{1}{\hat{\nu}} - 1) \begin{pmatrix} \hat{\mu}^* - \hat{\mu} \\ \hat{\nu}^* - \hat{\nu} \end{pmatrix} + \frac{n}{4}(\hat{\mu}^* - \hat{\mu}, \hat{\nu}^* - \hat{\nu}) \nabla^2 g(\hat{\mu}, \hat{\nu}) \begin{pmatrix} \hat{\mu}^* - \hat{\mu} \\ \hat{\nu}^* - \hat{\nu} \end{pmatrix} \\ &= z_1(z_1^* - z_1) - z_2(z_2^* - z_2) + \frac{1}{2}(z_1^* - z_1)^2 - \frac{1}{2}(z_2^* - z_2)^2.\end{aligned}\tag{S13}$$

35 Thus Equation S12 gives $P_1 = \mathbb{P}\{\Delta^* > 0|x\} = \mathbb{P}\{z_1^{*2} - z_2^{*2}|x\}$, as in Equation S2, and
36 Equation 15 follows.

37 Note that

$$P_1 = \mathbb{P}\{\Delta^* > 0|x\} = \mathbb{P}\{\Delta^* - \Delta > -\Delta|x\} = 1 - F_{\Delta^* - \Delta}(-\Delta),\tag{S14}$$

38 where $F(\cdot)$ is the CDF of $\Delta^* - \Delta$. We note that being the difference of two χ_d^2 variates, Δ has the
39 same distribution as $-\Delta$, which may and may not be the same distribution as that of $\Delta^* - \Delta$.
40 When $\Delta^* - \Delta$ and $-\Delta$ have the same distribution, $P_1 \sim \mathbb{U}(0, 1)$. Otherwise if $\Delta^* - \Delta$ and $-\Delta$ have
41 different distributions, P_1 does not have a uniform distribution. This follows from the fact that if
42 F is a smooth monotonic increasing function $\mathbb{R} \rightarrow (0, 1)$, and Y is a random variable over \mathbb{R} , then
43 $F(Y) \sim \mathbb{U}(0, 1)$ implies that F is the CDF of Y , because the CDF of Y is
44 $F_Y(y) = \mathbb{P}(Y < y) = \mathbb{P}(F(Y) < F(y)) = F(y)$.

45 According to DasGupta (2008, p. 475), if $T = \sqrt{n}(g(\bar{Y}) - g(\theta_*))$ and $\nabla g(\theta_*) = 0$, then T
46 and T^* have different distributions and bootstrap fails to estimate the CDF of T consistently. In
47 problem 4, $\nabla g(\mu_*, \nu_*) = 0$, so that $\Delta^* - \Delta$ has a different distribution from $\Delta - \Delta_*$, and P_1 is not
48 distributed as $\mathbb{U}(0, 1)$. In problems 1 and 2, $\Delta^* - \Delta$ and $\Delta - \Delta_*$ have the same distribution, so that
49 $P_1 \sim \mathbb{U}(0, 1)$.

Table S1. H_0 is true model while $H_k, k = 1, 2, 3, 4$ are candidate models.

Problem	Models	Bayesian	Bootstrap
(a1) Equally wrong and distinct, no parameters	Truth: $p = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$, $H_1 : p = (0.4, 0.3, 0.3)$, $H_2 : p = (0.3, 0.4, 0.3)$, $H_3 : p = (0.3, 0.3, 0.4)$.	Marginal 0-1 distributions: $\mathbb{P}(P_1 = 1) = 1/3$, $\mathbb{P}(P_2 = 1) = 1/3$, $\mathbb{P}(P_3 = 1) = 1/3$. $\mathbb{P}(P_i = 0) = 1 - \mathbb{P}(P_i = 1), i = 1, 2, 3$.	Joint density the same as for problem 5 (fig. 2). $E(P_1) = E(P_2) = E(P_3) = \frac{1}{3}$.
(a2) Equally wrong and distinct, no parameters	Truth: $p = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$, $H_1 : p = (0.4, 0.3, 0.3)$, $H_2 : p = (0.3, 0.4, 0.3)$, $H_3 : p = (a, b, b)$, with $a \approx 0.2708$ and $b \approx 0.3646$ to be the roots to $ab^2 = 0.0036, a + 2b = 1$.	Marginal 0-1 distributions: $\mathbb{P}(P_1 = 1) = \frac{5}{12}$, $\mathbb{P}(P_2 = 1) \approx 0.245$, $\mathbb{P}(P_3 = 1) \approx 0.338$, $\mathbb{P}(P_i = 0) = 1 - \mathbb{P}(P_i = 1), i = 1, 2, 3$.	Joint density in eq. S15. $\mathbb{E}(P_1) = \frac{5}{12}$, $\mathbb{E}(P_2) \approx 0.2455303 \dots$, $\mathbb{E}(P_3) \approx 0.337803 \dots$.
(b) Equally wrong and distinct, no parameters	Truth: $\mathbb{N}(0, 1)$, $H_1 : \mathbb{N}(\mu, 1), H_2 : \mathbb{N}(-\mu, 1)$, $H_3 : \mathbb{N}(0, 1/\tau_1), H_4 : \mathbb{N}(0, 1/\tau_2)$, $\mu \approx 0.7976, \tau_1 = 0.25, \tau_2 \approx 2.5866$, which satisfy $\log \frac{\tau_1}{\tau_2} = \tau_1 - \tau_2, \mu^2 = \tau_1 - \log \tau_1 - 1$.	Marginal 0-1 distributions: $\mathbb{P}(P_1 = 1) = \frac{1}{2\pi} (\arctan \frac{\sqrt{2\mu}}{\tau_2 - 1} + \arctan \frac{\sqrt{2\mu}}{1 - \tau_1})$, $\mathbb{P}(P_2 = 1) = \mathbb{P}(P_1 = 1)$, $\mathbb{P}(P_3 = 1) = \frac{1}{2\pi} (\pi - 2 \arctan \frac{\sqrt{2\mu}}{1 - \tau_1})$, $\mathbb{P}(P_4 = 1) = \frac{1}{2\pi} (\pi - 2 \arctan \frac{\sqrt{2\mu}}{\tau_2 - 1})$.	$\mathbb{E}(P_i), i = 1, 2, 3$, the same as left for Bayesian method.
(c) Equally wrong and indistinct, with parameters	Truth: $\mathbb{N}(0, 1)$, $H_1 : \mathbb{N}(\mu, 1/\tau_1), H_2 : \mathbb{N}(-\mu, 1/\tau_1)$, $H_3 : \mathbb{N}(\mu, 1/\tau_2), H_4 : \mathbb{N}(-\mu, 1/\tau_2)$, $\mu > 0, \tau_1 = 0.25, \tau_2 \approx 2.5866$.	$\mu \sim \text{Exp}(\lambda)$, $F_i(x) = \begin{cases} \frac{1}{2}, & x = 0, \\ \frac{1}{2} + \frac{1}{2} \Phi \left(\frac{\Phi^{-1}(x)}{\sqrt{\tau_i}} \right), & x > 0. \end{cases}$	The density of marginal distribution of $P_i, i = 1, 2, 3, 4$, is $f(p) = -\log(p)$. $\mathbb{E}(P_1) = \mathbb{E}(P_2) = \mathbb{E}(P_3) = \mathbb{E}(P_4) = \frac{1}{4}$.
(d) Equally right, with parameters	Truth: $\mathbb{N}(0, 1)$, $H_1 : \mathbb{N}(\mu, 1)$, $H_2 : \mathbb{N}(-\mu, 1)$, $H_3 : \mathbb{N}(0, 1/\tau), \mu > 0$.	$\mu \sim \text{Exp}(\lambda), \tau \sim \text{Exp}(\xi)$. $P_1 = \frac{\lambda \Phi(Z_1) \exp(\frac{1}{2} Z_1^2)}{\lambda \exp(\frac{1}{2} Z_1^2) + \sqrt{2\xi} e^{-\xi} \exp(\frac{1}{2} Z_2^2)}$, $P_2 = \frac{\lambda (1 - \Phi(Z_1)) \exp(\frac{1}{2} Z_1^2)}{\lambda \exp(\frac{1}{2} Z_1^2) + \sqrt{2\xi} e^{-\xi} \exp(\frac{1}{2} Z_2^2)}$, $P_3 = 1 - P_1 - P_2$, where $Z_1 \equiv \sqrt{\pi} \bar{x} \sim \mathbb{N}(0, 1)$, $Z_2 \equiv \sqrt{\frac{\pi}{2}} (s^2 - 1) \sim \mathbb{N}(0, 1)$.	Marginal densities are $f_1(p) = f_2(p) = -\log(p)$, $f_3(p) = -\log 1 - 2p $. $\mathbb{E}(P_1) = \mathbb{E}(P_2) = \frac{1}{4}$, $\mathbb{E}(P_3) = 0.5$.

50 Note.— For case a2, let $h = \frac{\log(0.3) - \log(b)}{\log(0.4) - \log(b)} \approx -2.102287$. The correlation matrix for $\Delta_{12}, \Delta_{13}, \Delta_{23}$
 51 is

$$\begin{pmatrix} 1 & \frac{\sqrt{3}}{2} & -\frac{h+2}{2\sqrt{h^2+h+1}} \\ & 1 & -\sqrt{\frac{3}{4}} \cdot \frac{h+1}{\sqrt{h^2+h+1}} \\ & & 1 \end{pmatrix}. \quad (\text{S15})$$

52 The expectations are

$$53 \quad \mathbb{E}(P_1) = \frac{1}{4} + \frac{1}{2\pi} \arcsin \rho \left(\frac{\Delta_{12}}{\sqrt{n}\sigma_{12}}, \frac{\Delta_{13}}{\sqrt{n}\sigma_{13}} \right) = \frac{5}{12},$$

$$54 \quad \mathbb{E}(P_2) = \frac{1}{4} + \frac{1}{2\pi} \arcsin \rho \left(-\frac{\Delta_{12}}{\sqrt{n}\sigma_{12}}, \frac{\Delta_{23}}{\sqrt{n}\sigma_{23}} \right) = 0.2455303\dots,$$

$$55 \quad \mathbb{E}(P_3) = \frac{1}{4} + \frac{1}{2\pi} \arcsin \rho \left(\frac{\Delta_{13}}{\sqrt{n}\sigma_{13}}, \frac{\Delta_{23}}{\sqrt{n}\sigma_{23}} \right) = 0.337803\dots.$$