

MULTISCALE BOOTSTRAP USING SCALEBOOT PACKAGE

HIDETOSHI SHIMODAIRA

1. INTRODUCTION

`scaleboot` is an add-on package for R. It is for calculating approximately unbiased (AU) p -values for a general problem from a set of multiscale bootstrap probabilities (BPs). Scaling is equivalent to changing the sample size of a data in bootstrap resampling. We compute BPs at several scales, from which a very accurate p -value is calculated (Shimodaira 2002). This multiscale bootstrap method has been implemented in `CONSEL` (Shimodaira and Hasegawa 2001) for phylogenetic inference and as the R add-on package `pvclust` (Suzuki and Shimodaira 2006) for hierarchical clustering. The point of the `scaleboot` package is to calculate an improved version of the AU p -value that is justified even for hypotheses with nonsmooth boundaries (Shimodaira 2008).

The basic usage of this package is illustrated in a simple example below. Then real applications in hierarchical clustering and phylogenetic inference are shown later.

For the use of `scaleboot`, Shimodaira (2008) may be referenced.

New update in 2018: Selective inference is now implemented. The selective inference version of the approximately unbiased p -value is denoted as SI in `scaleboot` and `pvclust`. The original approximately unbiased p -value, which is denoted as AU, is not designed for selective inference. The bootstrap probability (BP) is also not designed for selective inference. The theory and example of selective inference is described in Terada and Shimodaira (2017).

2. INSTALL

`scaleboot` is easily installed from CRAN online. Windows users can install the package by choosing “scaleboot” from the pull-down menu. Otherwise, run R on your computer and type

```
> install.packages("scaleboot")
```

You can also download the package file from the URL below and install it manually.

```
http://stat.sys.i.kyoto-u.ac.jp/prog/scaleboot/
```

3. SIMPLE EXAMPLE

3.1. Simulation Data. We first generate a simulation data.

```
> simdata <- function(n,y,sd) {  
+   m <- length(y)  
+   x <- matrix(rnorm(m*n,0,sd),m,n)  
+   t(x + (y - apply(x,1,mean)))  
+ }  
> X <- simdata(100,c(0,1,1,1,1,1,1,1,1,1),10)  
> round(X[1:3,],3)
```

This document is a part of the `scaleboot` package (Version 0.4-0 or newer). The source file is `usesb.Rnw`. I thank Paul A. Sheridan for his comments to improve the earlier version of the manuscript.

```

      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
[1,] -6.558 0.925 0.312 10.095 2.894 3.519 -4.529 7.589 -6.712 -2.012
[2,] -0.638 0.573 -0.915 8.625 2.958 0.040 -2.599 5.552 -7.598 24.689
[3,] -5.917 7.251 3.721 8.961 -6.419 -4.051 -5.913 2.753 -10.037 4.057
> y <- apply(X,2,mean)
> round(y,3)
[1] 0 1 1 1 1 1 1 1 1 1

```

The matrix $X = (x_{ij})$ above is of size $n \times m$ with $n = 100$, $m = 10$. We consider X as a data of sample size n , and rows $x_i = (x_{i1}, \dots, x_{im})$, $i = 1, \dots, n$, are observations of a random vector of m dimensions.

3.2. Null Hypothesis. Let μ be the unknown population mean of the row vectors. An estimate of μ is the sample average of the rows defined as $y = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$. Let $f(\mu)$ be a 0/1-valued (or false/true valued) function of μ . The null hypothesis we are going to test is represented as $f(\mu) = 1$. For example, $f(\mu) = 1$ if μ_1 is the largest among μ_1, \dots, μ_m , and $f(\mu) = 0$ otherwise. This $f(\mu)$ is implemented as `mc1(mu)` below.

```

> mc1 <- function(x) all(x[1] >= x[-1])
> mc1(y)
[1] FALSE

```

Although $f(y) = 0$ gives a rough idea whether $f(\mu) = 1$, we want to calculate a real number ranging between 0 and 1 which indicates the possibility of $f(\mu) = 1$. This is what `scaleboot` calculates as p -values.

3.3. Bootstrap Probabilities. A naive way to calculate a p -value is by bootstrap resampling. Let $X^* = (x_{ij}^*)$ be a bootstrap sample of X ; each row x_i^* is obtained by resampling with replacement from the rows x_1, \dots, x_n . Let n' be the size of the resampling so that X^* is a matrix of size $n' \times m$. The bootstrap replicate is $y^* = \bar{x}^* = \frac{1}{n'} \sum_{i=1}^{n'} x_i^*$. The following code generates an X^* with $n' = n$, and calculates $f(y^*)$. The resampling is made via a weight vector w ; w_i is the number of times that x_i is resampled in X^* .

```

> countw <- function(x,w,fn) {
+   y <- apply(w*x,2,sum)/sum(w)
+   fn(y)
+ }
> w <- as.vector(rmultinom(1,100,rep(1,100)))
> w
[1] 2 2 2 0 1 1 2 0 2 1 2 0 1 0 2 1 0 1 0 0 0 2 0 1 1 0 1 0 0 0 1 1 1 2 0 1 1
[38] 2 1 0 1 1 1 0 2 1 1 1 0 0 1 1 0 2 1 0 2 1 1 1 3 1 0 2 1 3 0 0 0 0 1 1 2 1
[75] 1 0 1 1 3 1 2 1 2 0 2 2 1 0 1 4 1 1 0 1 1 1 2 1 1 1
> countw(X,w,mc1)
[1] FALSE

```

Let B be the number of bootstrap samples we will generate, and y_1^*, \dots, y_B^* be the bootstrap replicates. Typically, $B = 10,000$. The BP is computed as $\sum_{i=1}^B f(y_i^*)/B$, where the ordinary BP uses $n' = n$. Since first introduced by Felsenstein (1985), it has been widely used as a p -value, but the bias is in fact rather large.

3.4. P -value Calculation. `scaleboot` calculates corrected p -values for improving BPs. First load the package by

```

> library(scaleboot)

```

Below, `sa` specifies the scales, and `nb` specifies B for each scale, so that $10,000 \times 13 = 130,000$ bootstrap samples are generated internally. It takes a few minutes on a pc.

```
> sa <- 9^seq(-1,1,length=13)
> nb <- 10000
> X.sb <- scaleboot(X,nb,sa,countw,mc1)
```

The main result (for $k = 3$) is shown by

```
> summary(X.sb) # k = 3 (default)
Raw Bootstrap Probability (scale=1) : 0.93 (0.10)
```

Hypothesis: null

Corrected P-values for Models (percent,Frequentist):

	k.3	sk.3	beta0	beta1	aic	weight
sing.3	37.68 (2.94)	42.97 (2.83)	1.64 (0.02)	0.68 (0.01)	-12.35	100.00
poly.3	17.23 (1.26)	22.84 (1.47)	1.67 (0.03)	0.59 (0.02)	15.02	
poly.2	6.34 (0.34)	9.70 (0.49)	1.92 (0.02)	0.39 (0.01)	177.96	
poly.1	0.07 (0.00)	0.13 (0.01)	3.21 (0.02)	0.00 (0.00)	3534.85	

Best Model: sing.3

Corrected P-values by the Best Model and by Akaike Weights Averaging:

	k.3	sk.3	beta0	beta1
best	37.68 (2.94)	42.97 (2.83)	1.64 (0.02)	0.68 (0.01)
average	37.68 (2.94)	42.97 (2.83)	1.64 (0.02)	0.68 (0.01)

The selective inference p -values (SI) are denoted as `sk.3`, and the non-selective AU p -values are denoted as `k.3`. We can also see additional results by specifying k ,

```
> summary(X.sb, k=1:3) # k = 1, 2, 3
```

Raw Bootstrap Probability (scale=1) : 0.93 (0.10)

Hypothesis: null

Corrected P-values for Models (percent,Frequentist):

	k.1	k.2	k.3	sk.1	sk.2	sk.3	beta0	beta1	aic	weight
sing.3	1.04 (0.04)	16.89 (0.82)	37.68 (2.94)	2.08 (0.09)	22.50 (0.97)	42.97 (2.83)	1.64 (0.02)	0.68 (0.01)	-12.35	100.00
poly.3	1.20 (0.05)	14.18 (0.93)	17.23 (1.26)	2.39 (0.10)	19.60 (1.15)	22.84 (1.47)	1.67 (0.03)	0.59 (0.02)	15.02	
poly.2	1.03 (0.05)	6.34 (0.34)	6.34 (0.34)	2.06 (0.09)	9.70 (0.49)	9.70 (0.49)	1.92 (0.02)	0.39 (0.01)	177.96	
poly.1	0.07 (0.00)	0.07 (0.00)	0.07 (0.00)	0.13 (0.01)	0.13 (0.01)	0.13 (0.01)	3.21 (0.02)	0.00 (0.00)	3534.85	

Best Model: sing.3

Corrected P-values by the Best Model and by Akaike Weights Averaging:

	k.1	k.2	k.3	sk.1	sk.2	sk.3	beta0	beta1
best	1.04 (0.04)	16.89 (0.82)	37.68 (2.94)	2.08 (0.09)	22.50 (0.97)	42.97 (2.83)	1.64 (0.02)	0.68 (0.01)
average	1.04 (0.04)	16.89 (0.82)	37.68 (2.94)	2.08 (0.09)	22.50 (0.97)	42.97 (2.83)	1.64 (0.02)	0.68 (0.01)

A class of AU p -values p_k indexed by $k = 1, 2, 3$, are calculated, and they are labelled as `k.1`, `k.2`, and `k.3`. The p -values are shown in percent, and the standard errors are given in parentheses. We should look at the row of **average** (the bottom line), and we can ignore the other rows. $p_1 \approx 1\%$ corresponds to the ordinary BP, and $p_2 \approx 18\%$ corresponds to the AU p -value of Shimodaira (2002). What we recommend to use here is $p_3 \approx 40\%$; this is the AU p -value of Shimodaira (2008). For this particular example, the common practice for calculating a p -value is to use the multiple comparisons method. If it is applied to y , the p -value is $p = 66\%$, which is rather close to p_3 in our example, whereas p_1 is obviously too small.

Internally, as explained in the next section, several models are fitted to the observed bootstrap probabilities. They are sorted in increasing order of AIC, and the best model is **sing.3**. The row of **sing.3** is duplicated in the row **best**, two lines from the bottom. The Akaike weights $\propto \exp(-AIC/2)$ are also computed for

models and the p -values are averaged using the weights so that we get the row of **average**, the bottom line. Typically, the two bottom lines are almost identical, and **average** is regarded as a smoothed version of **best** for small changes in data. I recommend to use **average** instead of **best**.

3.5. Internal Steps. We consider the following three steps (i)-(iii). Internally, the **scaleboot** function (i) performs the multiscale bootstrap, and (ii) estimates coefficients for candidate models. Then the **summary** method (iii) calculates the corrected p -values. These steps are explained below.

The results of steps (i) and (ii) are shown here.

```
> X.sb
Multiscale Bootstrap Probabilities (percent):
 1   2   3   4   5   6   7   8   9  10  11  12  13
0.00 0.01 0.05 0.12 0.29 0.68 0.93 1.57 2.21 2.77 3.40 3.94 4.69

Numbers of Bootstrap Replicates:
 1   2   3   4   5   6   7   8   9  10  11  12  13
10000 10000 10000 10000 10000 10000 10000 10000 10000 10000 10000 10000 10000

Scales (Sigma Squared):
 1   2   3   4   5   6   7 8   9  10  11  12  13
0.1111 0.1603 0.2309 0.3333 0.4808 0.6944 1 1.449 2.083 3.03 4.348 6.25 9.091

Coefficients:
      beta0      beta1      beta2
sing.3 1.1518 (0.1347) 1.1601 (0.1401) 0.8332 (0.1221)
poly.3 1.6337 (0.0284) 0.6569 (0.0210) -0.0318 (0.0024)
poly.2 1.9212 (0.0219) 0.3943 (0.0069)
poly.1 3.2056 (0.0182)

Model Fitting:
      rss      df pfit      aic
sing.3   7.65 10 0.6629 -12.35
poly.3  35.02 10 0.0001  15.02
poly.2 199.96 11 0.0000 177.96
poly.1 3558.85 12 0.0000 3534.85

Best Model:  sing.3
>
```

The results of (i) are the BPs for the 13 scales shown at first. Let α_{σ^2} denote the BP at scale σ^2 . Each BP is calculated from 10,000 bootstrap samples of size n' as the frequency of observing $f(y^*) = 1$. In **scaleboot**, n' is $\text{round}(n/\text{sa}[i])$, for $i = 1, \dots, 13$. Then, the scale is recalculated as $\sigma^2 = n/n'$ for taking account of the discreteness.

Step (ii) is performed by the **sbfit** function called internally from the **scaleboot** function for fitting parametric models to observed α_{σ^2} 's. By default, four models are considered as candidates; **poly.1**, **poly.2**, **poly.3**, and **sing.3**. Each of these models is denoted as $\psi(\sigma^2|\beta)$. Let $z_{\sigma^2} = \Phi^{-1}(1 - \alpha_{\sigma^2})$ be the bootstrap z -value at scale σ^2 , where $\Phi^{-1}(p) = \text{qnorm}(p)$. We work on $\sigma z_{\sigma^2}(y)$, which may be called a normalized bootstrap z -value. Considering σz_{σ^2} as a function of σ^2 , the coefficient vector β is estimated by fitting $\sigma z_{\sigma^2} = \psi(\sigma^2|\beta)$. Let $\hat{\beta}$ denote the estimated

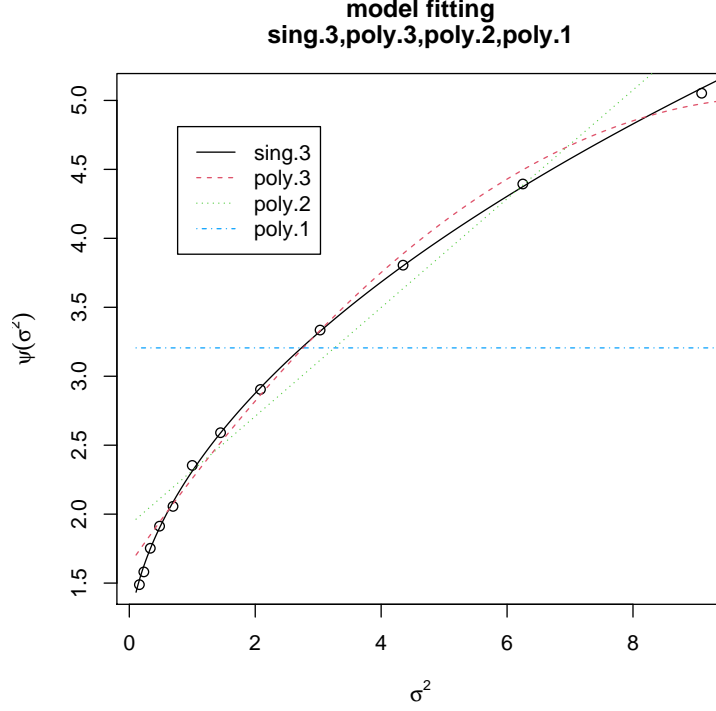


FIGURE 1. Model Fitting

value; the details of fitting $\hat{\beta}$ are explained later. We may choose the model which minimizes AIC value. The fitted curves are shown (Fig. 1) by plotting $\psi(\sigma^2|\hat{\beta})$ as

```
> plot(X.sb, legend="topleft")
```

The same plot but in other variables can be shown (Fig. 2) by, for example,

```
> plot(X.sb, xval="sigma", log="x", yval="pvalue", legend="topleft")
```

`poly.k` model is specified as a polynomial of σ^2 ; $\psi(\sigma^2|\beta) = \sum_{j=0}^{k-1} \beta_j \sigma^{2j}$ for $k \geq 1$. `sing.k` model is specified as $\psi(\sigma^2|\beta) = \beta_0 + \sum_{j=1}^{k-2} \beta_j \sigma^{2j} / (1 + \beta_{k-1}(\sigma - 1))$ for $k \geq 3$, where $0 \leq \beta_{k-1} \leq 1$. The number `k` for each model denotes the number of coefficients in β .

The details of model fitting are as follows. Let B_i and C_i be the number of replicates and the observed number of times that $f(y^*) = 1$, respectively, for the bootstrap resampling of scale σ_i^2 , $i = 1, \dots, S$. Since each C_i is binomially distributed, the log-likelihood is

$$\ell(\beta) = \sum_{i=1}^S \left\{ C_i \log \Phi(-\psi(\sigma_i^2|\beta)/\sigma_i) + (B_i - C_i) \log \Phi(\psi(\sigma_i^2|\beta)/\sigma_i) \right\},$$

where $\Phi(q) = \text{pnorm}(q)$. The estimate $\hat{\beta}$ is obtained by maximizing $\ell(\beta)$ numerically. The goodness of fit is measured by the difference of AIC values between the specified model and an unconstrained binomial model;

$$\text{AIC} = (-2\ell(\hat{\beta}) + 2k) - (-2\hat{\ell} + 2S),$$

where $\hat{\ell} = \sum_{i=1}^S (C_i \log(C_i/B_i) + (B_i - C_i) \log(1 - C_i/B_i))$.

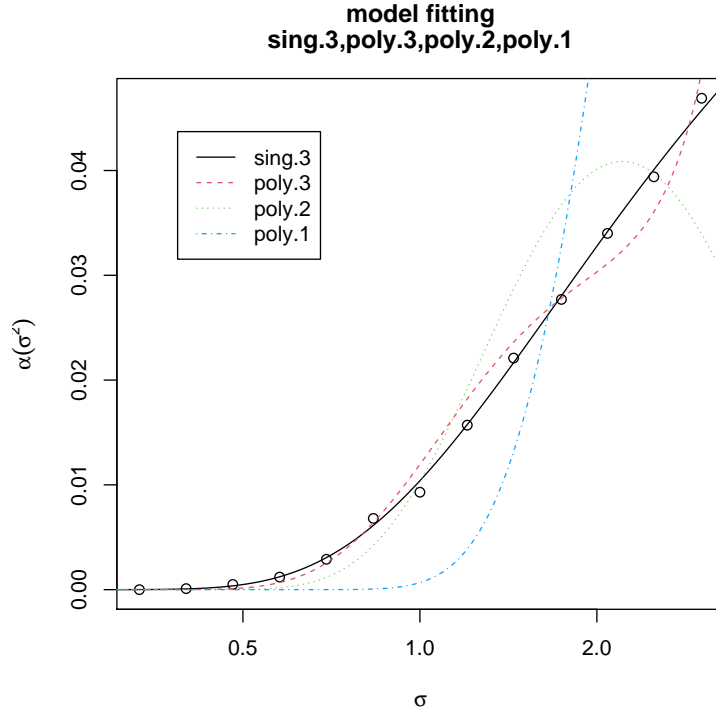


FIGURE 2. Model Fitting ($x = \log \sigma$, $y = \alpha_{\sigma^2}$)

Step (iii) is performed by the `summary` method as already mentioned. The first line shows the “raw” BP α_1 (the BP obtained from the ordinary bootstrap resampling). The main results are the corrected p -values, which follow next. For each model, we calculate q_k , $k = 1, 2, 3$, by

$$q_k = \sum_{j=0}^{k-1} \frac{(-1 - \sigma_0^2)^j}{j!} \frac{\partial^j \psi(\sigma^2 | \hat{\beta})}{\partial (\sigma^2)^j} \Big|_{\sigma_0^2}.$$

Then the corrected p -values are calculated by $p_k = 1 - \Phi(q_k)$. By default $\sigma_0^2 = 1$. The calculation of q_k is interpreted as extrapolation of σz_{σ^2} to $\sigma^2 = -1$ by using the first k terms of the Taylor series. According to the theory of Shimodaira (2008), the unbiased p -value is, if it exists, obtained by taking the limit $k \rightarrow \infty$. The extrapolated curves are shown (Fig. 3) by

```
> plot(summary(X.sb), legend="topleft")
```

4. HIERARCHICAL CLUSTERING

4.1. Pvclust Package. The `scaleboot` package includes an interface for the `pvclust` package (Suzuki and Shimodaira 2006). We use `pvclust` to calculate multiscale BPs for clusters by bootstrapping hierarchical clustering, from which we calculate an improved version of AU p -values using `scaleboot`. See `help(lung73)` for further details of the following example.

4.2. Using Pvclust. This example uses the `lung` dataset (Garber et al. 2001) included in `pvclust`. It is a DNA microarray data of 73 lung tissues (arrays) with 916 observations of genes. To draw dendrograms in terms of the arrays, we

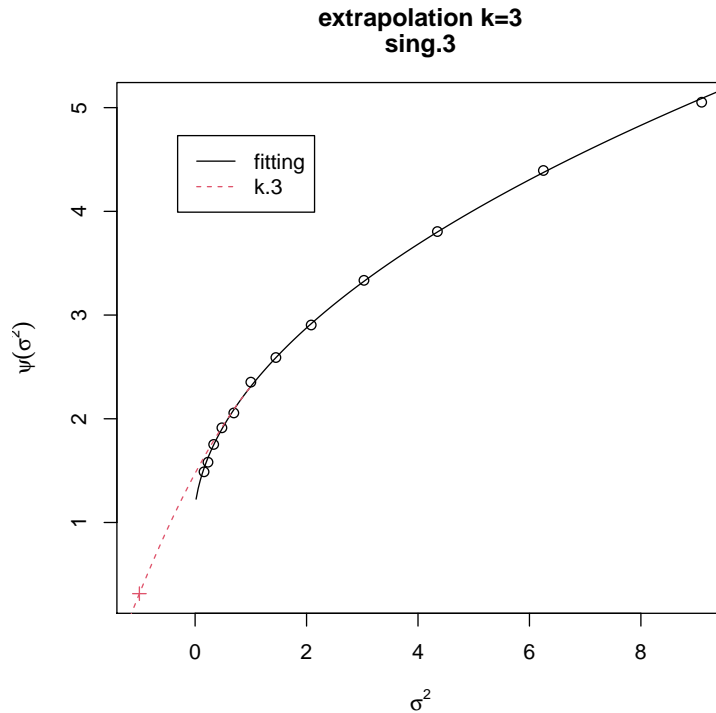


FIGURE 3. Extrapolation

resample genes in our analysis; this may be interpreted as assessing the uncertainty due to the variability of genes. The function `pvclust` first obtains a dendrogram by a hierarchical clustering method, and then calculates the multiscale BPs for each cluster of the dendrogram.

```
> library(pvclust)
> data(lung)
> sa <- 9^seq(-1,1,length=13)
> nb <- 10000
> lung73.pvclust <- pvclust(lung,r=1/sa,nboot=nb)
```

The above code may take a day, so it would be a good idea to run with `nb=1000` so that it would run 10 times faster. However, `nb=1000` should be used just for checking the program, and `nb=10,000` (at least) is recommended for publishing the results.

4.3. Model Fitting. We next apply the `sbfit` function of `scaleboot` to the multiscale BPs. For each cluster of the dendrogram, parametric models are fitted to the BPs.

```
> library(scaleboot)
> lung73.sb <- sbfit(lung73.pvclust)
```

4.4. Lung73 Dataset. The results of the previous two sections (`lung73.pvclust` and `lung73.sb`) are in fact stored in the `lung73` dataset of `scaleboot`. For users who want to try the examples, just type as follows.

```
> library(scaleboot)
> data(lung73)
```

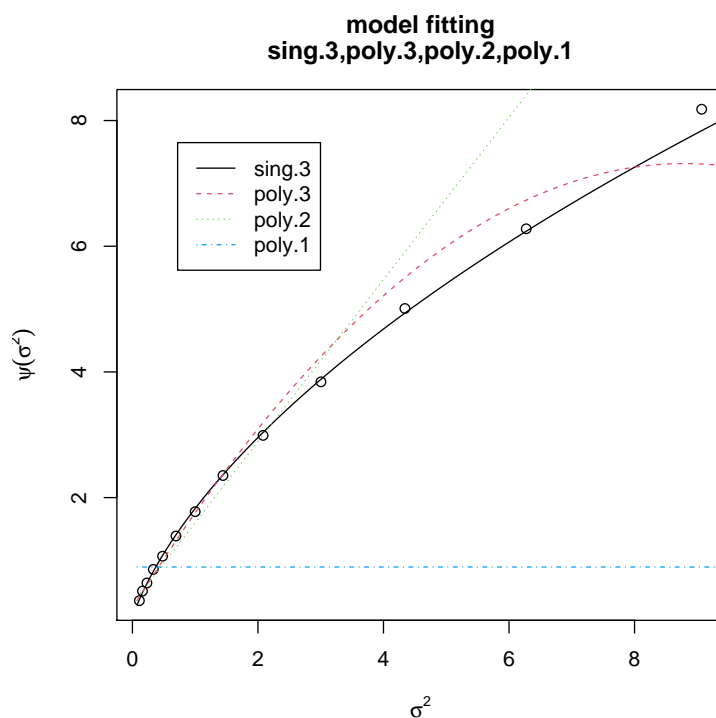



FIGURE 5. Model fitting for node 67

$k=1$, respectively, in tables generated by scaleboot for $k=3$. When you know the cluster before looking at the data (such as cluster id=31, the left most cluster in Fig. 4; it consists of normal cell and fetal cell, so no lung cancer at all), you may use AU. However, you should use SI for most of the clusters when they are identified by looking at the tree.

To calculate p_2 instead of p_3 , specify $k=2$,

```
> lung73.k2 <- sbpvcust(lung73.pvcust,lung73.sb, k=2)
```

4.6. Diagnostics of Fitting. The fitted curves are drawn by the plot method. For node 67, say, a plot with legend is obtained (Fig. 5) by

```
> plot(lung73.sb[[67]],legend="topleft")
```

All the calculated p -values for node 67 are given by

```
> summary(lung73.sb[[67]])
```

```
Raw Bootstrap Probability (scale=1) : 3.78 (0.19)
```

```
Hypothesis: null
```

```
Corrected P-values for Models (percent,Frequentist):
```

	k.3	sk.3	beta0	beta1	aic	weight
sing.3	95.00 (0.18)	97.45 (0.09)	0.53 (0.01)	1.29 (0.01)	34.35	100.00
poly.3	92.75 (0.28)	97.09 (0.14)	0.33 (0.00)	1.43 (0.01)	390.53	
poly.2	83.64 (0.32)	92.78 (0.19)	0.31 (0.00)	1.29 (0.01)	1319.11	
poly.1	18.53 (0.10)	37.05 (0.21)	0.90 (0.00)	0.00 (0.00)	53012.00	

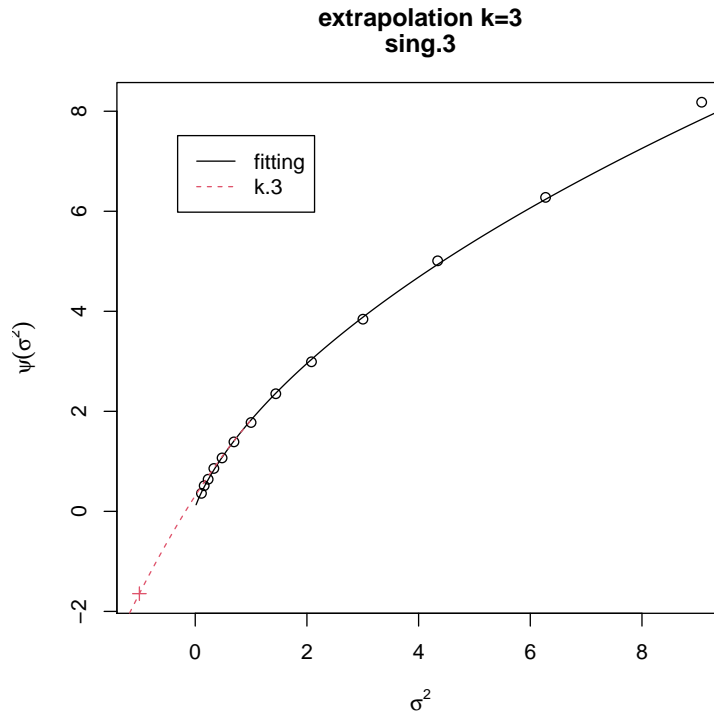


FIGURE 6. Extrapolation for node 67

Best Model: sing.3

Corrected P-values by the Best Model and by Akaike Weights Averaging:

	k.3	sk.3	beta0	beta1
best	95.00 (0.18)	97.45 (0.09)	0.53 (0.01)	1.29 (0.01)
average	95.00 (0.18)	97.45 (0.09)	0.53 (0.01)	1.29 (0.01)

The extrapolation using the best models (averaged by the Akaike weights) is shown (Fig. 6) by

```
> plot(summary(lung73.sb[[67]]), legend="topleft")
```

For a set of nodes, p -values are given by

```
> summary(lung73.sb[c(62,67,69,71)])
```

Corrected P-values by Akaike Weights Averaging (percent,Frequentist):

	raw	k.3	sk.3	beta0	beta1	hypothesis	model	weight
62	95.81 (0.20)	98.78 (0.07)	96.96 (0.16)	-2.00 (0.02)	0.25 (0.01)	alternative	poly.2	47.70
67	3.78 (0.19)	95.00 (0.18)	97.45 (0.09)	0.53 (0.01)	1.29 (0.01)	null	sing.3	100.00
69	30.51 (0.46)	75.04 (0.21)	9.00 (0.34)	-0.08 (0.00)	0.60 (0.00)	alternative	poly.2	57.35
71	25.63 (0.44)	86.02 (0.27)	26.18 (0.47)	-0.18 (0.00)	0.83 (0.01)	alternative	poly.3	100.00

Also plots are shown (Fig. 7) by

```
> plot(lung73.sb[c(62,67,69,71)])
```

5. PHYLOGENETIC INFERENCE

5.1. **CONSEL Software.** scaleboot has a front end for phylogenetic inference, and it may eventually replace the CONSEL software (Shimodaira and Hasegawa 2001)

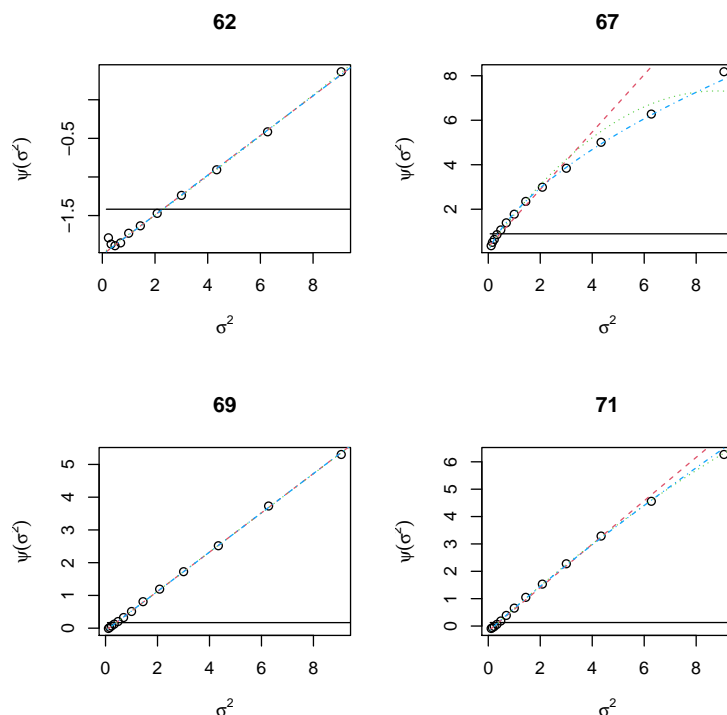


FIGURE 7. Model fitting for a set of nodes

for testing phylogenetic trees. Currently, `scaleboot` does not have a method for converting files obtained from other commonly used phylogenetic software packages, and so we must use `CONSEL` for this purpose before applying `scaleboot` to calculate an improved version of AU p -values for trees and edges. See `help(mam15)` for further details of the following example.

5.2. Mammal Dataset. We work on an example of phylogenetic analysis of six mammal species: *Homo sapiens* (human), *Phoca vitulina* (harbor seal), *Bos taurus* (cow), *Oryctolagus cuniculus* (rabbit), *Mus musculus* (mouse), *Didelphis virginiana* (opossum). The dataset was originally used in Shimodaira and Hasegawa (1999).

For Unix users, download `mam15-files.tgz`, and for Windows users download `mam15-files.zip`. The details of dataset files are as follows. `mam15.aa`: amino acid sequences ($n = 3414$) of mtDNA for the six mammals. `mam15.ass`: association vectors for edges and trees. `mam15.lnf`: site-wise log-likelihood values (output from PAML). `mam15.log`: detailed information for the associations. `mam15.mt`: site-wise log-likelihood values (output from seqmt). `mam15.tpl`: 15 tree topologies.

5.3. Likelihood Calculation of Trees. The main body of the dataset is the amino acid sequences (`mam15.aa`). We consider $m = 15$ tree topologies of the six mammals (`mam15.tpl`);

```
((Homsa, (Phovi, Bosta)), Orycu, (Musmu, Didvi)); t1
(Homsa, Orycu, ((Phovi, Bosta), (Musmu, Didvi))); t2
(Homsa, ((Phovi, Bosta), Orycu), (Musmu, Didvi)); t3
(Homsa, (Orycu, Musmu), ((Phovi, Bosta), Didvi)); t4
((Homsa, (Phovi, Bosta)), (Orycu, Musmu), Didvi); t5
```

```

(Homsa,((Phovi,Bosta),(Orycu,Musmu)),Didvi); t6
(Homsa,(((Phovi,Bosta),Orycu),Musmu),Didvi); t7
(((Homsa,(Phovi,Bosta)),Musmu),Orycu,Didvi); t8
(((Homsa,Musmu),(Phovi,Bosta)),Orycu,Didvi); t9
(Homsa,Orycu,(((Phovi,Bosta),Musmu),Didvi)); t10
(Homsa,(((Phovi,Bosta),Musmu),Orycu),Didvi); t11
((Homsa,((Phovi,Bosta),Musmu)),Orycu,Didvi); t12
(Homsa,Orycu,(((Phovi,Bosta),Didvi),Musmu)); t13
((Homsa,Musmu),Orycu,((Phovi,Bosta),Didvi)); t14
((Homsa,Musmu),((Phovi,Bosta),Orycu),Didvi); t15

```

The maximum likelihood estimates for these trees are calculated by PAML (Yang 1997). Let x_{ij} be the site-wise log-likelihood for sites $i = 1, \dots, n$, and trees $j = 1, \dots, m$. The log-likelihood of tree- j is $\sum_{i=1}^n x_{ij}$. A large n justifies the central limit theorem for $y = \bar{x}$, and allows us to resample x_{ij} directly without recalculation of the maximum likelihood estimates. The matrix $X = (x_{ij})$ is produced by PAML and stored in `mam15.lnf`. It is converted by CONSEL to a simpler format and stored in `mam15.mt`. The command is

```
seqmt --paml mam15.lnf
```

5.4. *P*-value Calculation for Trees. The AU *p*-values for trees are calculated simply by

```

> library(scaleboot)
> mam15.mt <- read.mt("mam15.mt")
> mam15.trees <- relltest(mam15.mt)
> summary(mam15.trees)

```

The `relltest` function above may take a half hour. The next section can be skipped if only tree selection is of interest.

5.5. *P*-value Calculation for Clusters. We can also calculate AU *p*-values for clusters (edges) of trees. We have to know, for each cluster, in which of the 15 trees it is included. The file `mam15.ass` has this information, which was generated using CONSEL by the command

```
treeass --outgroup 6 mam15.tpl > mam15.log
```

It also produces `mam15.log` for human readable information. A part of `mam15.log` is as follows.

```

# leaves: 6
6
  1 Homsa
  2 Phovi
  3 Bosta
  4 Orycu
  5 Musmu
  6 Didvi

# base edges: 10
10 6

    123456
  1 +++--- ;
  2 ++++-- ;
  3 +---+-- ;
  4 -++++-- ;
  5 ----+-- ;

```

```

6 +----- ;
7 -++++- ;
8 +++--- ;
9 +----- ;
10 -++++- ;

```

The clusters (edges) defined above are named e_1, \dots, e_{10} . For example, $e_1 = +++-- = (\text{Homsa}, \text{Phovi}, \text{Bosta})$.

The AU p -values for clusters as well as trees are calculated simply by

```

> library(scaleboot)
> mam15.mt <- read.mt("mam15.mt")
> mam15.ass <- read.ass("mam15.ass")
> mam15.relltest <- relltest(mam15.mt, ass=mam15.ass)
> summary(mam15.relltest)

```

5.6. Mam15 Dataset. The results of the previous sections (`mam15.mt`, `mam15.ass`, and `mam15.relltest`) are in fact stored in `mam15` dataset of `scaleboot`. For users who want to try the examples, just type as follows.

```

> library(scaleboot)
> data(mam15)

```

The results for trees are extracted by

```

> mam15.trees <- mam15.relltest[1:15]

```

We have used a cluster computer of 40 cpus for parallel computing using the `snow` package. The following code may take only 10 minutes, although we have used the number of resamples 10 times larger than the default value.

```

> library(snow)
> cl <- makeCluster(40)
> library(scaleboot)
> mam15.mt <- read.mt("mam15.mt")
> mam15.ass <- read.ass("mam15.ass")
> mam15.relltest <- relltest(mam15.mt, nb=100000, ass=mam15.ass)

```

5.7. Interpreting the Results. First we sort the results in increasing order of log-likelihood difference,

```

> stat <- attr(mam15.trees, "stat")
> o <- order(stat)
> mam15.trees <- mam15.trees[o]
> summary(mam15.trees, k=1:3)

```

Corrected P-values by Akaike Weights Averaging (percent, Frequentist):

	raw	k.1	k.2	k.3	sk.1	sk.2	sk.3	beta0	beta1	hypothesis	model	weight
t1	57.49 (0.16)	56.10 (0.04)	74.55 (0.05)	74.55 (0.06)	12.20 (0.08)	36.37 (0.09)	36.37 (0.09)	-0.41 (0.00)	0.25 (0.00)	alternative	poly.2	57.26
t3	31.88 (0.15)	30.37 (0.05)	46.39 (0.09)	45.21 (0.13)	60.73 (0.10)	79.46 (0.08)	78.61 (0.10)	0.30 (0.00)	0.21 (0.00)	null	poly.3	100.00
t2	3.67 (0.06)	3.65 (0.03)	13.10 (0.19)	16.60 (0.46)	7.31 (0.05)	20.75 (0.26)	24.66 (0.53)	1.46 (0.01)	0.34 (0.00)	null	poly.3	100.00
t5	1.32 (0.04)	1.35 (0.02)	7.83 (0.25)	10.13 (0.52)	2.70 (0.03)	11.97 (0.33)	14.62 (0.63)	1.81 (0.01)	0.40 (0.01)	null	poly.3	99.99
t6	3.21 (0.06)	3.16 (0.02)	12.44 (0.19)	13.60 (0.27)	6.32 (0.05)	19.51 (0.25)	20.82 (0.34)	1.51 (0.01)	0.35 (0.00)	null	poly.3	78.42
t7	0.54 (0.02)	0.52 (0.01)	3.55 (0.20)	4.52 (0.39)	1.04 (0.02)	5.48 (0.29)	6.66 (0.50)	2.18 (0.02)	0.38 (0.01)	null	poly.3	90.92
t4	1.54 (0.04)	1.49 (0.02)	10.28 (0.27)	14.59 (0.66)	2.97 (0.03)	15.23 (0.35)	19.98 (0.75)	1.72 (0.01)	0.45 (0.01)	null	poly.3	99.07
t15	0.08 (0.01)	0.07 (0.00)	0.96 (0.10)	1.17 (0.15)	0.15 (0.01)	1.45 (0.14)	1.71 (0.20)	2.76 (0.03)	0.42 (0.01)	null	poly.3	96.14
t8	0.00 (0.00)	0.00 (0.00)	0.05 (0.03)	0.10 (0.09)	0.00 (0.00)	0.08 (0.04)	0.14 (0.11)	3.73 (0.11)	0.42 (0.05)	null	poly.3	60.23
t14	0.24 (0.02)	0.22 (0.01)	2.87 (0.27)	5.18 (0.83)	0.45 (0.02)	4.21 (0.37)	6.91 (0.97)	2.37 (0.03)	0.47 (0.02)	null	poly.3	97.92
t13	0.02 (0.00)	0.01 (0.00)	0.38 (0.10)	0.81 (0.38)	0.02 (0.00)	0.55 (0.14)	1.04 (0.44)	3.19 (0.06)	0.49 (0.03)	null	poly.3	72.69
t9	0.00 (0.00)	0.00 (0.00)	0.09 (0.06)	0.32 (0.32)	0.00 (0.00)	0.13 (0.08)	0.39 (0.37)	3.67 (0.11)	0.52 (0.06)	null	poly.3	56.43
t11	0.00 (0.00)	0.00 (0.00)	0.02 (0.02)	0.05 (0.06)	0.00 (0.00)	0.03 (0.02)	0.07 (0.08)	4.05 (0.13)	0.49 (0.05)	null	poly.3	77.19
t10	0.00 (0.00)	0.00 (0.00)	0.00 (0.00)	0.00 (0.00)	0.00 (0.00)	0.00 (0.00)	0.00 (0.00)	5.51 (0.33)	0.37 (0.11)	null	poly.2	48.24
t12	0.00 (0.00)	0.00 (0.00)	0.00 (0.00)	0.00 (0.00)	0.00 (0.00)	0.00 (0.00)	0.00 (0.00)	5.45 (0.38)	0.45 (0.13)	null	poly.2	38.55

Next we look at the p -values. We confirm that p_1 (the second column, indicated as `k.1`) is almost the same as the raw BP (the first column); this should be so if the model fitting is good. Only two trees, i.e., `t1` and `t3`, have $p_1 > 0.05$. It is known that the bias of p_1 is large so that often leads to false positives for tree

selection. p_2 improves upon p_1 by correcting the bias. Six trees, i.e., $\mathbf{t1}$, $\mathbf{t3}$, $\mathbf{t2}$, $\mathbf{t5}$, $\mathbf{t6}$, and $\mathbf{t4}$, have $p_2 > 0.05$. p_3 improves upon p_2 even more, although the trees of $p_3 > 0.05$ are the same six trees in this example. In general, the accuracy of p_k (AU p -values, indicated as **k.k** in the table) increases for larger k , but it tends to have larger variation, so typically $k = 2$ or $k = 3$ would be used. Terada and Shimodaira (2017) introduced selective inference p -value, denoted as SI, or **sk.k** in the table. In the phylogenetic tree selection, this may only apply to the tree selected by the algorithm. In our case, the maximum likelihood tree ($\mathbf{t1}$), we can look at $p = 0.36$ (**sk.2**). By taking this tree as the alternative hypothesis, we can interpret $1 - p = 0.64$ as the selective p -value. So, we can reject the null ($\mathbf{t1}$ is not true) and claim that $\mathbf{t1}$ is the true tree when $p > 0.95$ at $\alpha = 0.05$. Unfortunately, $p = 0.36$ is too small and we cannot say anything in this table.

Next we look at p -values for clusters instead of trees.

```
> mam15.edges <- mam15.relltest[16:25] # 10 edges
> summary(mam15.edges,k=1:3)
```

Corrected P-values by Akaike Weights Averaging (percent,Frequentist):													
raw	k.1	k.2	k.3	sk.1	sk.2	sk.3	beta0	beta1	hypothesis	model	weight		
e1	58.82 (0.16)	58.07 (0.05)	71.77 (0.07)	72.09 (0.11)	16.14 (0.10)	33.76 (0.09)	34.08 (0.12)	-0.39 (0.00)	0.19 (0.00)	alternative	poly.3	100.00	
e2	93.04 (0.08)	93.06 (0.04)	95.59 (0.07)	95.65 (0.09)	86.12 (0.07)	90.31 (0.13)	90.41 (0.16)	-1.59 (0.00)	0.11 (0.00)	alternative	poly.3	60.86	
e3	3.68 (0.06)	3.63 (0.03)	12.62 (0.17)	17.96 (0.53)	7.26 (0.05)	20.11 (0.23)	26.02 (0.59)	1.47 (0.00)	0.33 (0.00)	null	sing.3	100.00	
e4	32.51 (0.15)	31.76 (0.05)	43.39 (0.09)	43.23 (0.12)	63.51 (0.10)	77.31 (0.08)	77.20 (0.10)	0.32 (0.00)	0.15 (0.00)	null	poly.3	98.52	
e5	6.07 (0.08)	5.91 (0.03)	7.38 (0.10)	7.12 (0.12)	11.81 (0.06)	14.12 (0.16)	13.77 (0.19)	1.51 (0.00)	0.06 (0.00)	null	poly.3	100.00	
e6	1.80 (0.04)	1.74 (0.02)	7.03 (0.18)	11.74 (0.68)	3.48 (0.04)	11.25 (0.25)	16.70 (0.78)	1.79 (0.01)	0.32 (0.01)	null	sing.3	100.00	
e7	3.75 (0.06)	3.73 (0.03)	9.17 (0.17)	10.62 (0.34)	7.45 (0.05)	15.56 (0.24)	17.31 (0.44)	1.56 (0.01)	0.23 (0.00)	null	sing.3	99.89	
e8	0.00 (0.00)	0.00 (0.00)	0.05 (0.02)	0.12 (0.10)	0.01 (0.00)	0.08 (0.03)	0.16 (0.12)	3.65 (0.08)	0.34 (0.04)	null	sing.3	60.70	
e9	0.32 (0.02)	0.30 (0.01)	1.46 (0.12)	2.28 (0.34)	0.60 (0.02)	2.39 (0.17)	3.41 (0.44)	2.47 (0.02)	0.28 (0.01)	null	poly.3	50.77	
e10	0.00 (0.00)	0.00 (0.00)	0.01 (0.01)	0.02 (0.02)	0.00 (0.00)	0.01 (0.01)	0.02 (0.02)	4.20 (0.12)	0.40 (0.04)	null	poly.3	83.98	

When we apply hypothesis testing to all the 10 clusters, we may use AU p -values. Let p be **k.2**. Then **e8**, **e9** and **e10** have $p < 0.05$, and **e5**, **e6**, **e7** also have $p < 0.10$. They may be rejected at either $\alpha = 0.05$ or 0.10 . However, people often want to look at only the clusters appearing in the top tree. In this table, only **e1** and **e2** are included in $\mathbf{t1}$. We should use SI p -values for testing each of these two clusters by taking it as the alternative hypothesis. Let p be **sk.2**. Then **e2** has $p > 0.90$, indicating that the null hypothesis (**e2** is not true) is rejected as $1 - p < 0.10$ at $\alpha = 0.10$.

Finally we examine model fitting. According to the AIC values, the fitting is good overall except for the top two trees; however note that the AIC values should be about 10 times smaller if the default value of $\text{nb}=10,000$ was used. The fitting curves for the top four trees are shown (Fig. 8) by

```
> plot(mam15.trees[1:4])
```

According to the plots, the fitting is rather good even for the top two trees.

REFERENCES

- [1] Felsenstein, J. (1985). Confidence limits on phylogenies: an approach using the bootstrap. *Evolution* **39** 783–791.
- [2] Garber, M. E. et al. (2001) Diversity of gene expression in adenocarcinoma of the lung. *Proceedings of the National Academy of Sciences* **98** 13784–13789 (dataset is available from http://genome-www.stanford.edu/lung_cancer/adeno/).
- [3] Shimodaira, H. (2002). An approximately unbiased test of phylogenetic tree selection. *Syst. Biol.* **51** 492–508.
- [4] Shimodaira, H. (2008) Testing Regions with Nonsmooth Boundaries via Multiscale Bootstrap. *Journal of Statistical Planning and Inference* **138** 1227–1241 (<http://dx.doi.org/10.1016/j.jspi.2007.04.001>).
- [5] Shimodaira, H. and Hasegawa, M. (2001) CONSEL: for assessing the confidence of phylogenetic tree selection. *Bioinformatics* **17** 1246–1247 (software is available from <http://stat.sys.i.kyoto-u.ac.jp/prog/consel/>).
- [6] Suzuki, R. and Shimodaira, H. (2006) pvclust: An R package for hierarchical clustering with p -values. *Bioinformatics* **22** 1540–1542 (software is available from <http://stat.sys.i.kyoto-u.ac.jp/prog/pvclust/>).

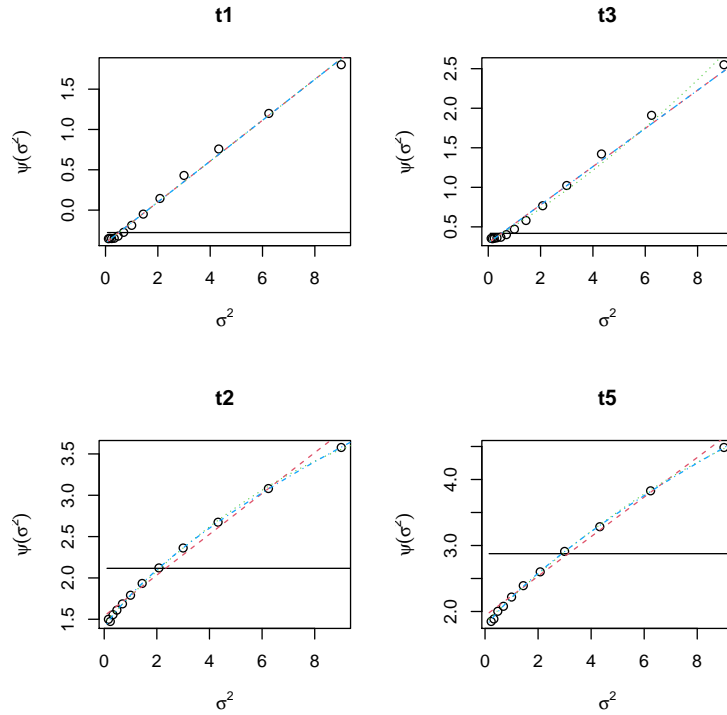


FIGURE 8. Model fitting for the top four trees

- [7] Terada, R. and Shimodaira, H. (2017) Selective inference for the problem of regions via multiscale bootstrap. *arXiv:1711.00949*.
- [8] Yang, Z. (1997) PAML: a program package for phylogenetic analysis by maximum likelihood. *Computer Applications in BioSciences* **13** 555–556 (software is available from <http://abacus.gene.ucl.ac.uk/software/paml.html>).

GRADUATE SCHOOL OF INFORMATICS, KYOTO UNIVERSITY, YOSHIDA HONMACHI, SAKYO-KU, KYOTO, 606-8501, JAPAN
 Email address: shimo@i.kyoto-u.ac.jp