

Growing phylogenetic trees with TreeLine

Erik S. Wright

July 16, 2024

Contents

1	Introduction	1
2	Performance Considerations	1
3	Growing a Phylogenetic Tree	2
4	Plotting Branch Support Values	5
5	Ancestral State Reconstruction	11
6	Exporting the Tree	13
7	Session Information	13

1 Introduction

This document describes how to grow phylogenetic trees using the `TreeLine` function in the DECIPHER package. `TreeLine` takes as input a set of aligned nucleotide or amino acid sequences and returns a phylogenetic tree (i.e., *dendrogram* object) as output. This vignette focuses on building maximum likelihood (ML) and maximum parsimony (MP) phylogenetic trees starting from sequences, but `TreeLine` can also be used to build additive trees from a distance matrix.

Why is the function called `TreeLine`? The goal of `TreeLine` is to find the most likely/parsimonious tree for a given sequence alignment. There are often many trees with nearly maximal likelihood/parsimony. Therefore, `TreeLine` seeks to find a tree as close as possible to the treeline, analogous to how no trees can grow above the treeline on a mountain.

Why use `TreeLine` versus other programs? The `TreeLine` function is designed to return an excellent phylogenetic tree with minimal user intervention. Many tree building programs have a large set of complex options for niche applications. In contrast, `TreeLine` simply builds a great tree when relying on its defaults. This vignette is intended to get you started and introduce additional options/functions that might be useful.

`TreeLine` uses multi-start optimization followed by hill-climbing to find the highest trees on the likelihood or parsimony landscapes. Since `TreeLine` is a stochastic optimizer, it optimizes many trees to prevent luck from influencing the final result. The main difference from most other approaches to tree optimization is that `TreeLine` heavily uses past trees to generate new trees as it optimizes. With any luck it'll find the treeline!

2 Performance Considerations

Finding a tree with very high likelihood/parsimony is no easy feat. `TreeLine` systematically optimizes hundreds to thousands of candidate trees before returning the best one. This takes time, but there are things you can do to make it

go faster.

- Only use the sequences you need: *TreeLine* scales approximately quadratically with the number of sequences. Hence, limiting the number of sequences is a worthwhile consideration. In particular, always eliminate redundant sequences, as shown below, and remove any sequences that are not necessary. This concern is shared for all tree building programs, and *TreeLine* is no exception.
- Choose a model: Automatic model selection is a useful feature, but frequently this time-consuming step can be skipped. For most modestly large sets of nucleotide sequences, the "GTR+G4" model will be automatically selected. Typical amino acid sequences will tend to pick the "LG+G4" or "WAG+G4" models, unless the sequences are from a particular origin (e.g., mitochondria). Pre-selecting a subset of the available `MODELS` and supplying this as the *model* argument can save considerable time.
- Set a timeout: The `maxTime` argument specifies the (approximate) maximum number of hours you are willing to let *TreeLine* run. If you are concerned about the code running too long then simply specify this argument.
- Compile with OpenMP support: Significant speed-ups can be achieved with multi-threading using OpenMP. See the "Getting Started DECIPHERing" vignette for how to do this on your computer platform. Then you only need to set the argument `processors=NULL` and *TreeLine* will use all available processors.
- Compile for SIMD support: *TreeLine* is configured to make use of SIMD operations, which are available on some processors. The easiest way to enable SIMD is to add "-O3 -march=native" to the end of `PKG_CFLAGS` in the "DECIPHER/src/MAKEVARS" text file. Then, after recompiling, there may be an automatic speed-up on systems with SIMD support. Note that enabling SIMD makes the compiled code non-portable, so the code always needs to be compiled on the hardware being used.

3 Growing a Phylogenetic Tree

TreeLine takes as input a multiple sequence alignment when constructing a maximum likelihood or maximum parsimony phylogenetic tree. Multiple sequence alignments can be constructed from a set of (unaligned) sequences using `AlignSeqs` or related functions. *TreeLine* will optimize trees for amino acid (i.e., `AAStringSet`) or nucleotide (i.e., `DNAStringSet` or `RNAStringSet`) sequences. Here, we are going to use a set of sequences that is included with DECIPHER. These sequences are from the internal transcribed spacer (ITS) between the 16S and 23S ribosomal RNA genes in several *Streptomyces* species.

```
> library(DECIPHER)
> # specify the path to your sequence file:
> fas <- "<<path to FASTA file>>"
> # OR find the example sequence file used in this tutorial:
> fas <- system.file("extdata", "Streptomyces_ITS_aligned.fas", package="DECIPHER")
> seqs <- readDNAStringSet(fas) # use readAAStringSet for amino acid sequences
> seqs # the aligned sequences
DNAStringSet object of length 88:
      width seq
[1] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTTTCCGAATGGGGAAACC supercont3.1 of S...
[2] 627 NNNNCACACCGCCCGTCA-CGTC...GGGGTTTCCGAATGGGGAAACC supercont3.1 of S...
[3] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTTTCCGAATGGGGAAACC supercont1.1 of S...
[4] 627 CGTACACACCGCCCGTCA-CGTC...GGGGTTTCCGAATGGGGAAACC supercont1.1 of S...
[5] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTTTCCGAATGGGGAAACC supercont1.1 of S...
...
[84] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTTTCCGAATGGGGAAACC gi|297189896|ref|...
[85] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTGTCCGAATGGGGAAACC gi|224581106|ref|...
```

```
[86] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTGTCCGAATGGGGAAACC gi|224581106|ref|...
[87] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTGTCCGAATGGGGAAACC gi|224581106|ref|...
[88] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTTTCCGAATGGGGAAACC gi|224581108|ref|...
```

Many of these sequences are redundant or from the same genome. We can de-replicate the sequences to accelerate tree building:

```
> seqs <- unique(seqs) # remove duplicated sequences
> ns <- gsub("^. *Streptomyces( subsp\\\. | sp\\\. | | sp_) ([^ ]+).*$", "\\2", names(seqs))
> names(seqs) <- ns # name by species (or any other preferred names)
> seqs <- seqs[!duplicated(ns)] # remove redundant sequences from the same species
> seqs
DNASet object of length 19:
      width seq
[1] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTTTCCGAATGGGGAAACC albus
[2] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTTTCCGAATGGGGAAACC clavuligerus
[3] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTGTCCGAATGGGGAAACC ghanaensis
[4] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTTTCCGAATGGGGAAACC griseoflavus
[5] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTGTCCGAATGGGGAAACC lividans
...
[15] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTGTCCGAATGGGGAAACC cattleya
[16] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTTTCCGAATGGGGAAACC bingchengensis
[17] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTTTCCGAATGGGGAAACC avermitilis
[18] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTGTCCGAATGGGGAAACC C
[19] 627 TGTACACACCGCCCGTCA-CGTC...GGGGTGTCCGAATGGGGAAACC Tu6071
```

Now, it's time to try our luck at finding the most likely tree. Here, we will set a stringent time limit (0.01 hours) to make this example faster, although longer time limits (e.g., 24 hours) are advised because setting very short time limits leaves the result partly up to luck.

Note that *TreeLine* automatically selects a substitution model based on Akaike information criterion (by default). It is possible to specify specific model(s) (e.g., `model="GTR+G4"`) to limit the possible selections and test your luck with fewer models.

Also, since *TreeLine* is a stochastic optimizer, it is critical to always set the random number seed for reproducibility. You can pick any lucky number, and if you ever wonder how much you pushed your luck, you can try running again from a different random number seed to see how much the result came down to luck of the draw.

Note that setting a time limit, as done below with *maxTime*, negates the purpose of setting a seed – never set a time limit if reproducibility is desired or you'll have no such luck.

```
> set.seed(123) # set the random number seed
> tree <- TreeLine(seqs,
  method="ML",
  model="GTR+G4",
  reconstruct=TRUE,
  maxTime=0.01)
```

Optimizing model parameters:

GTR+G4 -ln(L) = 4368, AICc = 8832, BIC = 9017

PHASE 1 OF 3: INITIAL TREES

```
1/3. Optimizing initial tree #1 of 10 to 100:
-ln(L) = 4366.5 (-0.043%), 2 Climbs
1/3. Optimizing initial tree #2 of 10 to 100:
-ln(L) = 4365.2 (-0.029%), 7 Climbs
1/3. Optimizing initial tree #3 of 11 to 100:
-ln(L) = 4363.8 (-0.034%), 6 Climbs
```

PHASE 2 OF 3: REGROW GENERATION 1 OF 10 TO 20

```
2/3. Optimizing regrown tree #1 of 10 to 100:
-ln(L) = 4366.5 (+0.062%), 1 Climb
2/3. Optimizing regrown tree #2 of 10 to 100:
-ln(L) = 4377.0 (+0.303%), 3 Climbs
2/3. Optimizing regrown tree #3 of 10 to 100:
-ln(L) = 4366.5 (+0.063%), 2 Climbs
2/3. Optimizing regrown tree #4 of 10 to 100:
-ln(L) = 4367.5 (+0.085%), 2 Climbs
2/3. Optimizing regrown tree #5 of 10 to 100:
-ln(L) = 4363.8 (~0.000%), 2 Climbs
2/3. Optimizing regrown tree #6 of 10 to 100:
-ln(L) = 4363.8 (~0.000%), 2 Climbs
2/3. Optimizing regrown tree #7 of 10 to 100:
-ln(L) = 4366.5 (+0.062%), 2 Climbs
2/3. Optimizing regrown tree #8 of 10 to 100:
-ln(L) = 4364.4 (+0.015%), 3 Climbs
```

PHASE 3 OF 3: SHAKEN TREES

Grafting 1 tree to the best tree:

-ln(L) = 4363.8 (0.000%), 0 Grafts of 0

```
3/3. Optimizing shaken tree #1 of 3 to 1000:
-ln(L) = 4364.1 (+0.009%), 8 Climbs
3/3. Optimizing shaken tree #2 of 3 to 1000:
-ln(L) = 4492.6 (+2.868%), 9 Climbs
3/3. Optimizing shaken tree #3 of 3 to 1000:
-ln(L) = 4385.5 (+0.496%), 9 Climbs
```

Grafting 3 trees to the best tree:

-ln(L) = 4363.8 (0.000%), 0 Grafts of 7

Model parameters:

Frequency(A) = 0.175

Frequency(C) = 0.244

Frequency(G) = 0.347

Frequency(T) = 0.234

Rate A <-> C = 0.711

4 Plotting Branch Support Values

Maybe it was just beginner's luck, but we already have a reasonable looking starting tree! TreeLine automatically returns a variety of information about the tree that can be accessed with the `attributes` and `attr` functions:

```
> #attributes(tree) # view all attributes
> attr(tree, "members") # number of leaves below this (root) node
[1] 19
> attr(tree, "height") # height of the node (in this case, the midpoint root)
[1] 2.24726
> attr(tree, "state") # ancestral state reconstruction (if reconstruct=TRUE)
[1] "-----CACCGCCCGTCA-CGTCACGAAAGTCGGTAACACCCGAAGCCGGTGGCCCAACCCCTTG-GGGAGGGAGCTGTCGAA
> attr(tree, "siteLnLs") # LnL for every alignment column (site)
  [1] -2.032764 -1.575973 -2.032764 -2.356720 -1.949814 -2.356720
  [7] -2.093367 -2.507668 -2.093367 -2.093367 -1.707977 -2.093367
 [13] -2.093367 -2.093367 -1.707977 -2.180270 -2.093367 -2.507668
 [19]  0.000000 -2.093367 -1.707977 -2.180270 -2.093367 -2.507668
 [25] -4.588630 -1.707977 -2.507668 -2.507668 -2.507668 -1.707977
 [31] -2.180270 -2.093367 -1.707977 -1.707977 -2.180270 -2.507668
 [37] -5.896149 -4.588630 -5.896149 -2.093367 -2.093367 -2.093367
 [43] -1.707977 -2.507668 -2.507668 -1.707977 -2.093367 -2.093367
 [49] -5.072818 -4.451366 -2.180270 -1.707977 -1.707977 -2.093367
 [55] -2.093367 -2.093367 -2.507668 -2.507668 -2.093367 -2.093367
 [61] -2.093367 -12.830845 -8.042482 -16.165681 -13.554422 -6.702884
 [67] -1.707977 -1.707977 -1.707977 -2.507668 -1.707977 -1.707977
 [73] -1.707977 -2.507668 -1.707977 -8.438611 -11.715708 -1.707977
 [79] -2.180270 -2.093367 -1.707977 -2.507668 -2.507668 -1.707977
 [85] -1.707977 -2.180270 -1.707977 -1.707977 -1.707977 -2.507668
 [91] -2.093367 -13.681408 -12.549894 -1.707977 -2.093367 -1.707977
 [97] -2.507668 -2.180270 -2.180270 -1.707977 -1.707977 -1.707977
[103] -2.507668 -2.093367 -1.707977 -2.507668 -2.507668 -1.707977
[109] -2.180270 -2.093367 -1.707977 -2.180270 -2.507668 -2.507668
[115] -2.093367 -2.507668 -2.507668 -1.707977 -1.707977 -2.180270
[121] -2.507668 -1.707977 -2.093367 -2.093367 -1.707977 -2.180270
[127] -2.507668 -2.093367 -2.093367 -1.707977 -1.707977 -2.507668
[133] -2.507668 -1.707977 -1.707977 -2.180270 -1.707977 -2.093367
[139] -1.707977 -1.707977 -2.093367 -2.180270 -1.707977 -1.707977
[145] -2.507668 -2.180270 -2.093367 -2.507668 -2.093367 -2.093367
[151] -2.180270 -2.093367 -2.093367 -2.180270 -2.180270 -2.180270
[157] -2.093367 -2.180270 -2.507668 -2.507668 -1.707977 -1.707977
[163] -2.507668 -1.707977 -2.093367 -2.507668 -14.367073 -19.521379
[169] -5.560379 -12.145905 -12.350502 -13.471003 -17.207677 -13.711659
[175] -11.241465 -17.541286 -7.625300 -9.041922 -8.363745 -15.844685
[181] -22.474908 -19.872629 -7.612833 -12.223259 -11.490907 -10.707143
[187] -8.805690 -12.913811 -9.786175 -11.416556 -14.211128 -11.096251
[193] -6.830006 -7.274107 -6.051166 -10.120789 -9.357439 -8.549835
[199] -14.291896 -11.040292 -13.904419 -11.948841 -8.848408 -2.075651
[205] -7.777453 -6.320082 -17.565072 -15.623646 -16.439435 -9.298684
[211] -15.568964 -10.233744 -13.848452 -14.583496 -6.493554 -2.093367
[217] -15.379594 -10.520892 -13.369443 -16.847630 -1.707977 -15.404783
```

[223]	-17.483819	-15.631559	-16.599544	-14.711668	-18.243885	-12.005012
[229]	-15.777375	-9.692514	-1.707977	-4.873445	-2.180270	-15.550874
[235]	-2.426490	-4.451366	-2.093367	-5.748567	-2.093367	-6.093714
[241]	-11.410647	-1.707977	-5.072818	-4.451366	-4.732445	-1.707977
[247]	-5.072818	-4.795675	-4.795675	-4.588630	-1.707977	-2.180270
[253]	-2.180270	-1.707977	-6.093714	-16.839045	-5.121046	-4.795675
[259]	-15.779709	-4.732445	-7.535339	-4.873445	-1.707977	-7.132231
[265]	-11.121890	-15.297997	-18.846885	-19.616929	-13.354832	-21.669123
[271]	-20.353888	-22.004579	-24.656106	-23.372423	-20.963712	-18.861277
[277]	-22.895856	-25.045341	-23.464786	-26.983363	-20.860484	-23.947126
[283]	-17.827201	-11.577203	-18.320183	-12.906152	-18.850334	-21.728886
[289]	-4.795675	-1.707977	-2.180270	-5.896149	-5.765911	-2.180270
[295]	-6.998524	-5.162073	-11.379001	-13.619443	-13.829721	-6.232887
[301]	-5.362810	-7.250735	-10.644021	-10.503949	-8.242438	-1.707977
[307]	-13.602805	-6.892924	-9.728638	-1.707977	-10.635730	-9.234879
[313]	-8.990857	-25.511400	-17.238950	-15.545015	-1.410746	-1.452262
[319]	-1.452262	-7.385522	-21.401475	-24.083503	-22.172397	-22.173370
[325]	-26.561666	-24.039128	-23.263644	-21.756078	-23.152328	-20.020659
[331]	-12.647556	-3.747955	-19.056771	-23.682666	-23.970821	-18.339273
[337]	-20.982809	-23.168351	-17.541364	-12.854535	-17.462450	-9.146776
[343]	-16.916378	-5.162073	-8.830206	-1.707977	-1.707977	-5.765911
[349]	-8.453544	-2.093367	-7.944282	-12.498180	-4.732445	-1.707977
[355]	-5.748567	-5.121046	-1.707977	-1.410746	-1.707977	-1.707977
[361]	-5.748567	-13.552286	-12.109769	-2.093367	-2.180270	-1.707977
[367]	-4.795675	-10.563691	-1.707977	-1.707977	-23.988771	-12.778582
[373]	-7.305050	-1.707977	-8.998280	-17.336955	-4.463027	-16.642707
[379]	-17.851399	-4.504062	-1.678386	-1.452262	-21.749354	-21.178944
[385]	-23.033443	-8.418358	-20.138963	-12.733074	-22.864803	-14.564883
[391]	-20.507854	-16.274203	-15.185030	-17.117237	-11.143287	-17.458938
[397]	-24.059734	-16.154167	-1.938177	-1.410746	-1.452262	-1.452262
[403]	-9.942793	-19.544643	-12.469499	-9.337496	-4.588630	-5.162073
[409]	-1.707977	-12.638978	-10.524040	-2.093367	-6.784770	-4.588630
[415]	-18.083042	-8.785957	-7.717436	-17.123321	-33.384638	-2.507668
[421]	0.000000	-9.446735	-9.164492	-13.957892	-21.434650	-29.084628
[427]	-21.421239	-20.005131	-22.026654	-17.361296	-18.230091	-18.052038
[433]	-21.420034	-30.259792	-14.884946	-24.799359	-23.296423	-24.253300
[439]	-22.795162	-14.337522	-19.284270	-16.291081	-1.707977	-7.894649
[445]	-7.622297	-15.888112	-18.948617	-4.873445	-1.707977	-1.707977
[451]	-7.894649	-1.707977	-5.635654	-16.037483	-6.643286	-1.575973
[457]	-1.575973	-2.032764	-4.588630	-1.707977	-5.121046	-2.180270
[463]	-1.707977	-12.147370	-2.180270	-2.180270	-1.707977	-2.507668
[469]	-1.707977	-2.507668	-2.507668	-2.093367	-9.253764	-10.884957
[475]	-4.588630	-2.507668	-12.042784	-2.507668	-1.707977	-2.180270
[481]	-1.707977	-1.707977	-2.507668	-4.588630	-1.707977	-2.093367
[487]	-8.658988	-2.507668	-1.707977	-2.093367	-2.507668	-2.180270
[493]	-2.093367	-2.180270	-1.452262	-1.452262	-1.707977	-2.180270
[499]	-1.707977	-1.707977	-4.588630	-2.093367	-2.507668	-2.507668
[505]	-1.707977	-2.180270	-2.180270	-5.121046	-2.180270	-2.180270
[511]	-2.507668	-2.507668	-1.707977	-1.707977	-1.707977	-2.093367
[517]	-4.451366	-2.093367	-2.507668	-4.588630	-1.707977	-1.707977
[523]	-2.180270	-1.707977	-1.707977	-2.507668	-2.180270	-1.707977

```

[529] -4.588630 -2.093367 -2.180270 -2.180270 -1.707977 -1.707977
[535] -10.838946 -5.896149 -4.588630 -2.093367 -2.507668 -1.707977
[541] -1.707977 -2.507668 -4.795675 -2.093367 -2.093367 -1.707977
[547] -2.507668 -2.180270 -1.707977 -2.507668 -2.507668 -1.707977
[553] -1.707977 -2.507668 -2.093367 -1.707977 -2.180270 -1.707977
[559] -6.493554 -1.707977 -2.507668 -1.707977 -1.707977 -2.093367
[565] -4.588630 -10.308886 -2.093367 -1.707977 -2.507668 -2.180270
[571] -2.507668 -4.873445 -17.960848 -2.093367 -2.093367 -4.588630
[577] -2.093367 -1.707977 -1.707977 -1.707977 -1.707977 -2.507668
[583] -1.707977 -17.720916 -10.928175 -1.707977 -7.019384 -2.093367
[589] -2.507668 -2.507668 -2.093367 -6.770807 -10.308886 -11.366141
[595] -1.707977 -2.093367 -2.180270 -11.347053 -2.180270 -1.707977
[601] -2.507668 -2.180270 -2.093367 -2.093367 -1.707977 -4.451366
[607] -1.707977 -1.707977 -4.451366 -2.180270 -18.353619 -2.180270
[613] -2.093367 -2.093367 -1.707977 -2.507668 -2.507668 -2.180270
[619] -1.707977 -1.707977 -1.707977 -1.707977 -2.507668 -2.507668
[625] -2.507668 -2.093367 -2.093367
> attr(tree, "score") # best score (in this case, the -LnL)
[1] 4363.752
> attr(tree, "model") # either the specified or automatically select transition model
[1] "GTR+G4"
> attr(tree, "parameters") # the free model parameters (or NA if unoptimized)
      FreqA      FreqC      FreqG      FreqT      FreqI      A/G      C/T      A/C
0.1748244 0.2436608 0.3471494      NA      NA 3.0921698 2.9448983 0.7107601
      A/T      C/G      Indels      alpha
1.1107077 0.5897411      NA 0.1936143
> attr(tree, "midpoint") # center of the edge (for plotting)
[1] 9.61499

```

The tree is (virtually) rooted at its midpoint by default. For maximum likelihood trees, all internal nodes include aBayes branch support values [1]. These are given as probabilities that can be used in plotting on top of each edge. We can also italicize the leaf labels (species names).

```

> plot(dendrapply(tree,
  function(x) {
    s <- attr(x, "probability") # choose "probability" (aBayes) or "support"
    if (!is.null(s) && !is.na(s)) {
      s <- formatC(as.numeric(s), digits=2, format="f")
      attr(x, "edgetext") <- paste(s, "\n")
    }
    attr(x, "edgePar") <- list(p.col=NA, p.lwd=1e-5, t.col="#CC55AA", t.cex=
    if (is.leaf(x))
      attr(x, "nodePar") <- list(lab.font=3, pch=NA)
    x
  })),
  horiz=TRUE,
  yaxt='n')
> # add a scale bar (placed manually)
> arrows(0, 0, 0.4, 0, code=3, angle=90, len=0.05, xpd=TRUE)
> text(0.2, 0, "0.4 subs./site", pos=3, xpd=TRUE)

```

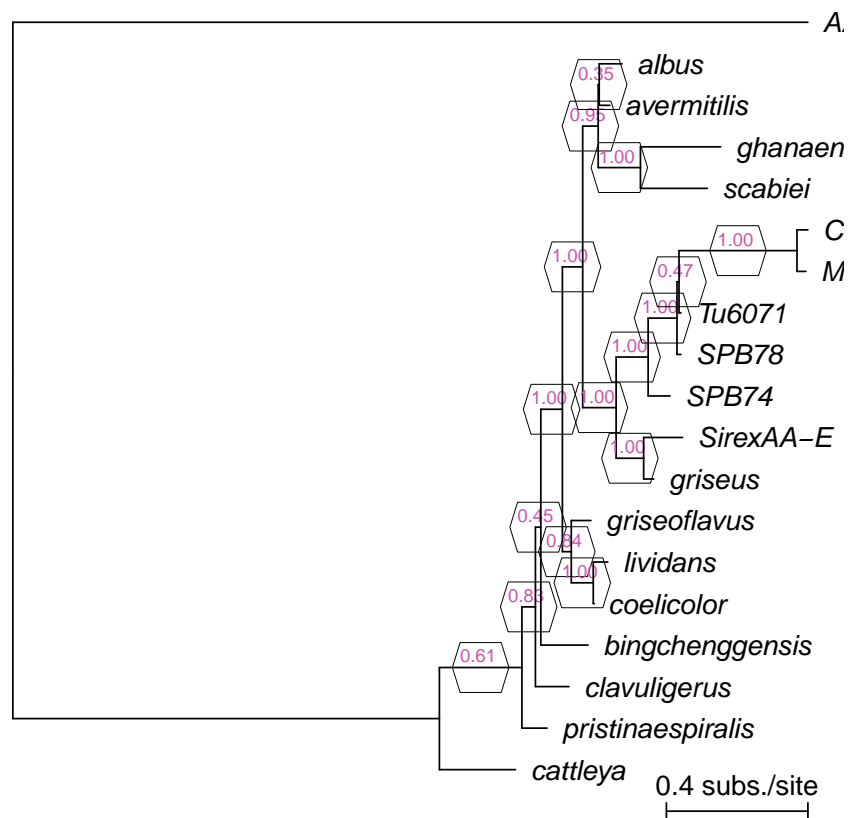


Figure 2: Tree with (aBayes) support probabilities at each internal node.

We lucked out because maximum likelihood and maximum parsimony trees both provide branch supports in the form of the fraction of optimized trees that contained a given partition (branch). These are accessible from the “support” attribute. As luck would have it, support values and (aBayes) probabilities are correlated, but support tends to be more conservative.

```

> getSupports <- function(x) {
  if (is.leaf(x)) {
    NULL
  } else if (is.null(attr(x, "support"))) {
    rbind(getSupports(x[[1]]), getSupports(x[[2]]))
  } else {
    rbind(cbind(attr(x, "support"),
                 attr(x, "probability"),
                 attr(x, "members")),
           getSupports(x[[1]]), getSupports(x[[2]]))
  }
}
> support <- getSupports(tree)
> plot(support[, 1],
       support[, 2],
       cex=log10(support[, 3]),
       xlab="Support", ylab="aBayes probability", asp=1)
> abline(a=0, b=1, lty=2) # line of identity (y=x)

```

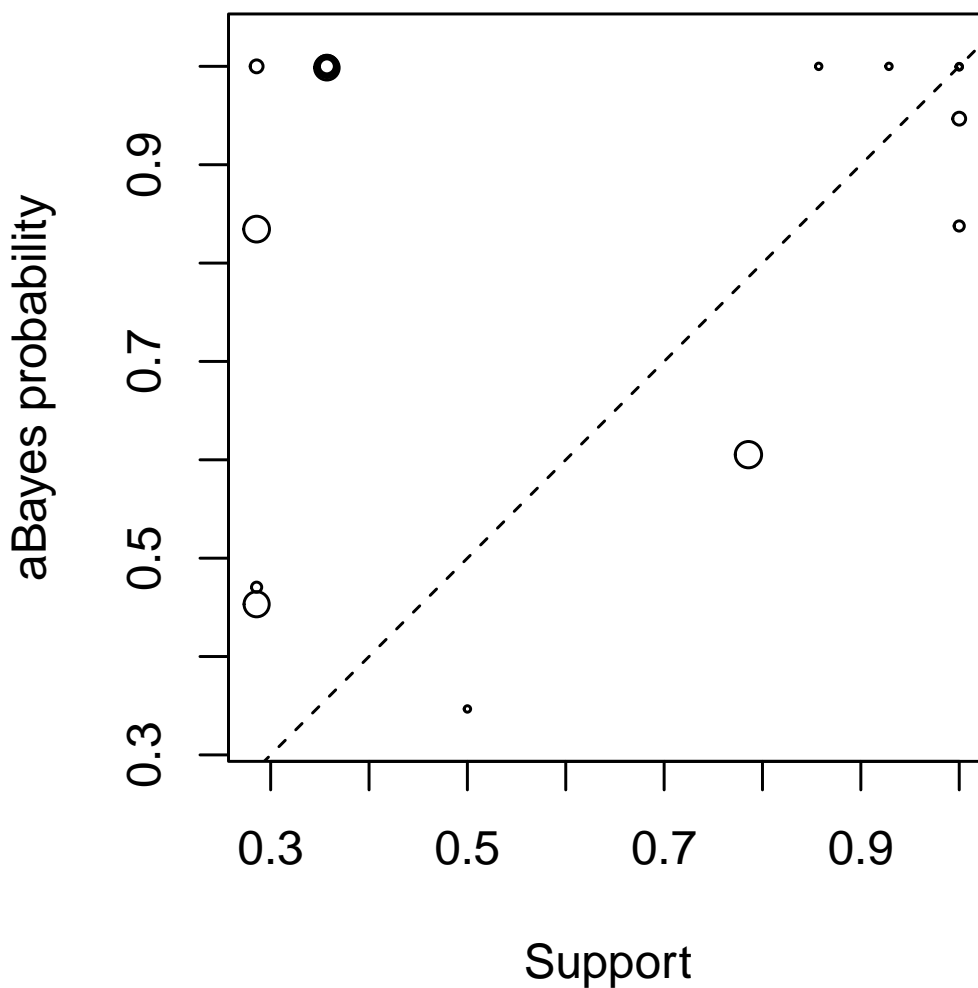


Figure 3: Comparison of aBayes probabilities and branch support values.

5 Ancestral State Reconstruction

We’re in luck—one of the advantages of maximum likelihood and maximum parsimony tree building methods is that they automatically predict states at each internal node on the tree [2]. This feature is enabled when *reconstruct* is set to `TRUE`. These character states can be used by the function `MapCharacters` to determine state transitions along each edge of the tree. This information enables us to plot the total number of substitutions occurring along each edge. The state transitions can be accessed along each edge by querying a new “change” attribute.

```

> new_tree <- MapCharacters(tree, labelEdges=TRUE)
> plot(new_tree, edgePar=list(p.col=NA, p.lwd=1e-5, t.col="#55CC99", t.cex=0.7))
> attr(new_tree[[1]], "change") # state changes on first branch left of (virtual) root
[1] "G168A" "G171T" "G176C" "G182C" "G184T" "G185T" "G199A" "G208C" "G214A"
[10] "G224T" "G225C" "G227C" "G229C" "G272T" "G276T" "G277T" "G279T" "G280T"
[19] "G283T" "G287C" "G288T" "G302C" "G303T" "G304C" "G314C" "G316T" "G321T"
[28] "G324T" "G325C" "G326C" "G327T" "G328T" "G330A" "G333C" "G338T" "G371A"
[37] "G379A" "G385T" "G387C" "G389C" "G391C" "G394A" "G396A" "G397T" "G419A"
[46] "G435A" "G447T" "G584C"

```

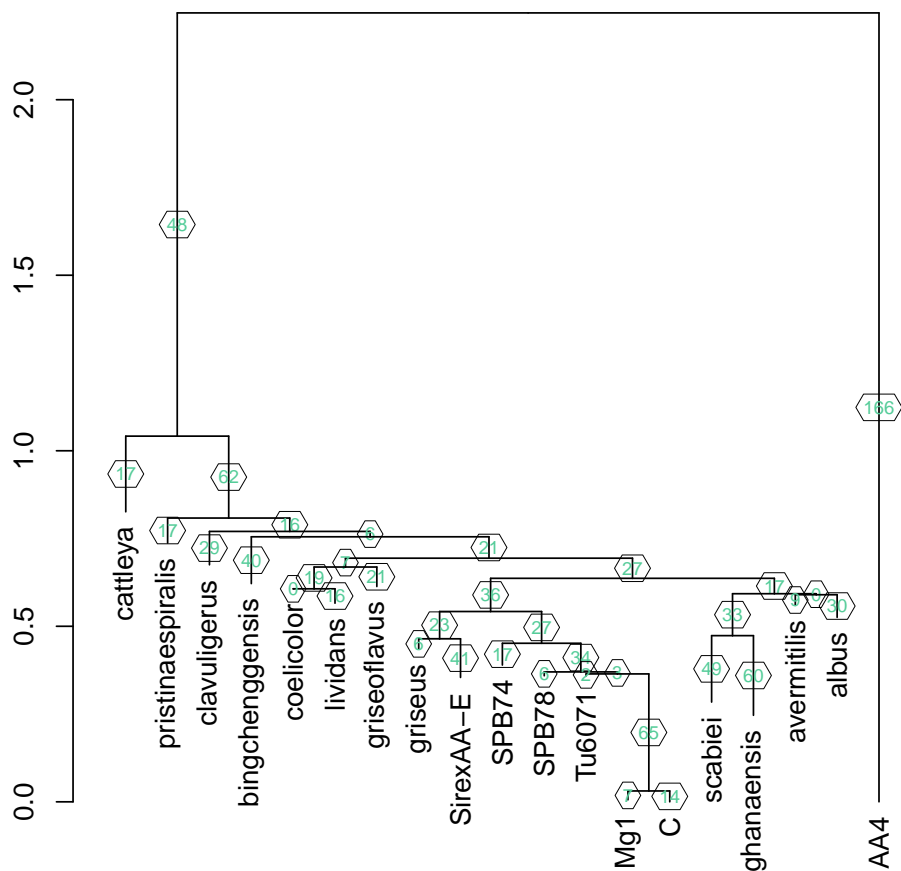


Figure 4: Edges labeled with the number of state transitions.

6 Exporting the Tree

We’ve had a run of good luck with this tree, so we’d better save it before our luck runs out! The functions `ReadDendrogram` and `WriteDendrogram` will import and export trees in Newick file format. If we leave the *file* argument blank then it will print the output to the console for our viewing:

```
> WriteDendrogram(tree, file="")
(('cattleya':0.2156051,('pristinaespiralis':0.07112583,('clavuligerus':0.09506737,('bing
```

To keep up our lucky streak, we should probably include any model parameters in the output along with the tree. Luckily, Newick format supports square brackets (i.e., “[]”) for comments, which we can append to the end of the file for good luck:

```
> params <- attr(tree, "parameters")
> cat("[", paste(names(params), params, sep="=", collapse=","), "]", sep="", append=TRUE
[FreqA=0.174824373876447,FreqC=0.243660806026823,FreqG=0.347149403352744,FreqT=NA,FreqI=
```

7 Session Information

All of the output in this vignette was produced under the following conditions:

- R version 4.4.1 (2024-06-14 ucrt), x86_64-w64-mingw32
- Running under: Windows Server 2022 x64 (build 20348)
- Matrix products: default
- Base packages: base, datasets, grDevices, graphics, methods, stats, stats4, utils
- Other packages: BiocGenerics 0.51.0, Biostrings 2.73.1, DECIPHER 3.1.4, GenomeInfoDb 1.41.1, IRanges 2.39.1, S4Vectors 0.43.1, XVector 0.45.0
- Loaded via a namespace (and not attached): DBI 1.2.3, GenomeInfoDbData 1.2.12, KernSmooth 2.23-24, R6 2.5.1, UCSC.utils 1.1.0, compiler 4.4.1, crayon 1.5.3, httr 1.4.7, jsonlite 1.8.8, tools 4.4.1, zlibbioc 1.51.1

References

- [1] Anisimova, M., Gil, M., Dufayard, J., Dessimoz, C., & Gascuel, O. Survey of branch support methods demonstrates accuracy, power, and robustness of fast likelihood-based approximation schemes. *Syst Biol.*, 60(5), 685-699.
- [2] Joy, J., Liang, R., McCloskey, R., Nguyen, T., & Poon, A. Ancestral Reconstruction. *PLoS Comp. Biol.*, 12(7), e1004763.