Tutorial using the \mathbf{Q} software

Genetic data analysis using \mathbf{Q} : introduction to phylogenetics

Thibaut Jombart

Abstract

This tutorial aims to illustrate the basics of phylogenetic reconstruction using **@**. Different kinds of phylogenetic approaches are introduced, namely distance-based, maximum parsimony, and maximum likelihood methods. We also illustrate how to assess the quality of phylogenetic trees using simple approaches. Methods are illustrated using a toy dataset of seasonal influenza isolates sampled in the US from 1993 to 2008.

Contents

1	Introduction									
	1.1	Required packages	3							
	1.2	The data	3							
2	Distance-based phylogenies									
	2.1	Computing genetic distances	5							
	2.2	Building trees	7							
	2.3	Plotting trees	9							
	2.4	Assessing the quality of a phylogeny	12							
3	Max	ximum parsimony phylogenies	18							
	3.1	Introduction	18							
	3.2	Implementation	18							

4	Maximum likelihood phylogenies									
	4.1	Introduction	20							
	4.2	Sorting out the data	21							
	4.3	Getting a ML tree	23							

1 Introduction

1.1 Required packages

This tutorial requires a working version of \mathbb{R} [5] greater than or equal to 2.12.1. It uses the following packages: *stats* implements basic hierarchical clustering routines, *ade4* [1] and *adegenet* [2] are here used essentially for their graphics, *ape* [4] is the core package for phylogenetics, and *phangorn* [6] implements parsimony and likelihood based methods. Make sure that the dependencies are installed as well when installing the packages:

```
> install.packages("adegenet", dep = TRUE)
> install.packages("phangorn", dep = TRUE)
```

Then load the packages using:

```
> library(stats)
> library(ade4)
> library(ape)
> library(adegenet)
> library(phangorn)
```

Some graphical functions used in this tutorial are also only part of the devel version of *adegenet*, and may not be present in the installed version of the package. To make sure these functions are available, source the patch online:

```
> source("http://adegenet.r-forge.r-project.org/files/patches/auxil.R")
```

or simply install the devel version using:

```
> install.packages("adegenet", repos = "http://r-forge.r-project.org")
```

1.2 The data

The data used in this tutorial are DNA sequences of seasonal influenza (H3N2) downloaded from Genbank (http://www.ncbi.nlm.nih.gov/genbank/). Alignments have been realized beforehand using standard tools (Clustalw2 for basic alignment and Jalview for refining the results). We selected 80 isolates genotyped for the hemagglutinin (HA) segment sampled in the US from 1993 to 2008. The dataset consists in two files: i) usflu.fasta, a file containing aligned DNA sequences and ii) usflu.annot.csv, a comma-separated file containing useful annotations of the sequences. In the following, we assume that both these files are stored in a data directory.

To read the DNA sequences into R, we use $\verb"read.dna"$ from the *ape* package:

```
> dna <- read.dna(file = "data/usflu.fasta", format = "fasta")
> dna
```

80 DNA sequences in binary format stored in a matrix. All sequences of same length: 1701 Labels: CY013200 CY013781 CY012128 CY013613 CY012160 CY012272 ... Base composition: a c g t 0.335 0.200 0.225 0.239

> class(dna)

```
[1] "DNAbin"
```

Sequences are stored as DNAbin objects, an efficient representation of DNA/RNA sequences which use bytes (as opposed to character strings) to code nucleotides:

```
> object.size(as.character(dna))/object.size(dna)
```

```
7.71879054549557 bytes
```

For instance, the first 10 nucleotides of the first 5 isolates:

```
> as.character(dna)[1:5, 1:10]
```

	[,1]	[,2]	[,3]	[,4]	[,5]	[,6]	[,7]	[,8]	[,9]	[,10]
CY013200	"a"	"t"	"ġ"	"a"	"a"	"g"	"a"	"c"	"t"	"a"
CY013781	"a"	"t"	"g"	"a"	"a"	"g"	"a"	"c"	"t"	"a"
CY012128	"a"	"t"	"g"	"a"	"a"	"g"	"a"	"c"	"t"	"a"
CY013613	"a"	"t"	"g"	"a"	"a"	"g"	"a"	"c"	"t"	"a"
CY012160	"a"	"t"	"g"	"a"	"a"	"ğ"	"a"	"c"	"t"	"a"

are actually coded as raw bytes:

> unclass(dna)[1:5, 1:10]

	[,1]	[,2]	[,3]	[,4]	[,5]	[,6]	[,7]	[,8]	[,9]	[,10]
CY013200	88	18	48	88	88	48	88	28	18	88
CY013781	88	18	48	88	88	48	88	28	18	88
CY012128	88	18	48	88	88	48	88	28	18	88
CY013613	88	18	48	88	88	48	88	28	18	88
CY012160	88	18	48	88	88	48	88	28	18	88

> typeof(unclass(dna)[1:5, 1:10])

[1] "raw"

The annotation file is read in R using the standard procedure:

```
> annot <- read.csv("data/usflu.annot.csv", header = TRUE, row.names = 1)
> head(annot)
```

accession year misc (A/New York/783/1993(H3N2)) CY013200 1993 2 1993 (A/New York/802/1993(H3N2)) CY013781 3 (A/New York/758/1993(H3N2)) CY012128 1993 4 CY013613 1993 (A/New York/766/1993(H3N2)) (A/New York/762/1993(H3N2)) 5 CY012160 1993 6 CY012272 1994 (A/New York/729/1994(H3N2))

accession contains the Genbank accession numbers, which are unique sequence identifiers; year is a year of collection of the isolates; misc contains other possibly useful information. Before going further, we check that isolates are identical in both files (accession number are used as labels for the sequences):

Good! The data we will analyse are 80 isolates (5 per year) typed for the same 1701 nucleotides.

2 Distance-based phylogenies

Distance-based phylogenetic reconstruction consits in i) computing pairwise genetic distances between individuals (here, isolates), ii) representing these distances using a tree, and iii) evaluating the relevance of this representation.

2.1 Computing genetic distances

We first compute genetic distances using *ape*'s dist.dna, which proposes no less than 15 different genetic distances (see ?dist.dna for details). Here, we use Tamura and Nei 1993's model [7] which allows for different rates of transitions and transversions, heterogeneous base frequencies, and betweensite variation of the substitution rate.

> D <- dist.dna(dna, model = "TN93")
> class(D)

[1] "dist"

> length(D)

[1] 3160

D is an object of class dist which contains the distances between every pairs of sequences.

Now that genetic distances between isolates have been computed, we need to visualize this information. There are n(n-1)/2 distances for n sequences, and most of the time summarising this information is not entirely trivial. The simplest approach is plotting directly the matrix of pairwise distances:

```
> temp <- as.data.frame(as.matrix(D))
> table.paint(temp, cleg = 0, clabel.row = 0.5, clabel.col = 0.5)
```



Darker shades of grey represent greater distances. Note that to use image to produce similar plots, data need to be transformed first; for instance:

```
> temp <- t(as.matrix(D))
> temp <- temp[, ncol(temp):1]
> par(mar = c(1, 5, 5, 1))
> image(x = 1:80, y = 1:80, temp, col = rev(heat.colors(100)),
+ xaxt = "n", yaxt = "n", xlab = "", ylab = "")
> axis(side = 2, at = 1:80, lab = rownames(dna), las = 2, cex.axis = 0.5)
> axis(side = 3, at = 1:80, lab = rownames(dna), las = 3, cex.axis = 0.5)
```



(see image.plot in the package *fields* for similar plots with a legend).

Since the data are roughly ordered by year, we can already see some genetic structure appearing, but this is admittedly not the most satisfying or informative approach, and tells us little about the evolutionary relationships between our isolates.

2.2 Building trees

We use trees to get a better representation of the genetic distances between individuals. It is important, however, to bear in mind that the obtained trees are not necessarily efficient representations of the original distances, and information can —and likely will— be lost in the process.

A wide array of algorithms for constructing trees from a distance matrix are available in \mathbb{Q} , including:

- * nj (ape package): the classical Neighbor-Joining algorithm.
- * bionj (ape): an improved version of Neighbor-Joining.
- * fastme.bal and fastme.ols (*ape*): minimum evolution algorithms.
- * hclust (*stats*): classical hierarchical clustering algorithms including single linkage, complete linkage, UPGMA, and others.

Here, we go for the standard:

```
> tre <- nj(D)
> class(tre)
[1] "phylo"
> tre
Phylogenetic tree with 80 tips and 78 internal nodes.
Tip labels:
        CY013200, CY013781, CY012128, CY013613, CY012160, CY012272, ...
Unrooted; includes branch lengths.
> plot(tre, cex = 0.6)
```

```
> title("A simple NJ tree")
```



A simple NJ tree

Trees created in the package *ape* are instances of the class phylo. See **?read.tree** for a description of this class.

2.3 Plotting trees

The plotting method offers many possibilities for plotting trees; see ?plot.phylo for more details. Functions such as tiplabels, nodelabels, edgelabels and axisPhylo can also be useful to annotate trees. For instance, we may simply represent years using different colors (red=ancient; blue=recent):

```
> plot(tre, show.tip = FALSE)
> title("Unrooted NJ tree")
> myPal <- colorRampPalette(c("red", "yellow", "green", "blue"))
> tiplabels(annot$year, bg = num2col(annot$year, col.pal = myPal),
+ cex = 0.5)
> temp <- pretty(1993:2008, 5)
> legend("bottomright", fill = num2col(temp, col.pal = myPal),
+ leg = temp, ncol = 2)
```



Unrooted NJ tree

This illustrates a common mistake when interpreting phylogenetic trees. In the above figures, we tend to assume that the left-side of the phylogeny is 'ancestral', while the right-side is 'recent'. This is wrong —as suggested by the colors— unless the phylogeny is actually rooted, i.e. some external taxa has been used to define what is the most 'ancient' split in the tree. The present tree is not rooted, and should be better represented as such:

```
> is.rooted(tre)
```

[1] FALSE

Unrooted NJ tree



In the present case, a sensible rooting would be any of the most ancient isolates (from 1993). We can take the first one:

```
> head(annot)
```

```
accession year
                                                 misc
   CY013200 1993 (A/New York/783/1993(H3N2))
CY013781 1993 (A/New York/802/1993(H3N2))
CY012128 1993 (A/New York/758/1993(H3N2))
1
2
3
   CY013613 1993 (A/New York/766/1993(H3N2))
4
   CY012160 1993 (A/New York/762/1993(H3N2))
CY012272 1994 (A/New York/729/1994(H3N2))
5
6
> tre2 <- root(tre, out = 1)
and plot the result:
> plot(tre2, show.tip = FALSE, edge.width = 2)
  >
>
+
> axisPhylo()
> temp <- pretty(1993:2008, 5)
> legend("topright", fill = transp(num2col(temp, col.pal = myPal),
       0.7), leg = temp, ncol = 2)
+
```



Rooted NJ tree

Now, the horizontal axis can globally be interpreted as temporal evolution; however, it is not uncommon that isolates from consecutive years cluster together, suggesting that the turnover of strains from one season to another is somehow smooth.

2.4 Assessing the quality of a phylogeny

Many genetic distances and hierarchical clustering algorithms can be used to build trees; not all of them are appropriate for a given dataset. Genetic distances rely on hypotheses about the evolution of DNA sequences which should be taken into account. For instance, the mere proportion of differing nucleotides between sequences (model='raw' in dist.dna) is easy to interprete, but only makes sense if all substitutions are equally frequent. In practice, simple yet flexible models such as that of Tamura and Nei (1993, [7]) are probably fair choices.

Once one has chosen an appropriate genetic distance and built a tree using this distance, an essential yet most often overlooked question is whether this tree actually is a good representation of the original distance matrix. This is easily investigated using simple biplots and correlation indices. The function **cophenetic** is used to compute distances between the tips of the tree. Note that more distances are available in the *adephylo* package (see **distTips** function).

```
> x <- as.vector(D)
> y <- as.vector(as.dist(cophenetic(tre2)))
> plot(x, y, xlab = "original distance", ylab = "distance in the tree",
+ main = "Is NJ appropriate?", pch = 20, col = transp("black",
+ 0.1), cex = 3)
> abline(lm(y ~ x), col = "red")
> cor(x, y)^2
```

```
[1] 0.9975154
```



Is NJ appropriate?

As it turns out, our Neighbor-Joining tree (tre2) is a very good representation of the chosen genetic distances. Things would have been different had we chosen, for instance, UPGMA:

```
> tre3 <- as.phylo(hclust(D, method = "average"))
> y <- as.vector(as.dist(cophenetic(tre3)))
> plot(x, y, xlab = "original distance", ylab = "distance in the tree",
+ main = "Is UPGMA appropriate?", pch = 20, col = transp("black",
+ 0.1), cex = 3)
> abline(lm(y ~ x), col = "red")
> cor(x, y)^2
```

[1] 0.7393009



Is UPGMA appropriate?

In this case, UPGMA is a poor choice. Why is this? A first explanation is that UPGMA forces ultrametry (all the tips are equidistant to the root):

> plot(tre3, cex = 0.5)
> title("UPGMA tree")

UPGMA tree



The underlying assumption is that all lineages have undergone the same amount of evolution, which is obviously not the case in seasonal influenza sampled over 16 years.

Another validation of phylogenetic trees, much more commonly used, in bootstrap. Bootstrapping a phylogeny consists in sampling the nucleotides with replacement, rebuilding the phylogeny, and checking if the original nodes are present in the bootstrapped trees. In practice, this procedure is iterated a large number of times (e.g. 100, 1000), depending on how computer-intensive the phylogenetic reconstruction is. The underlying idea is to assess the variability in the obtained topology which results from conducting the analyses on a random sample the genome. Note that the assumption that the analysed sequences represent a random sample of the genome is often dubious. For instance, this is not the case in our toy dataset, since HA segment has a different rate of evolution and experiences different selective pressures from other segments of the influenza genome. We nonetheless illustrate the procedure, implemented by **boot.phylo**:

```
> myBoots <- boot.phylo(tre2, dna, function(e) root(nj(dist.dna(e,</pre>
      model = "TN93")), 1))
+
> myBoots
           23
25
                                                                    42
71
                                          97 100 100
 [1] 100
               55
                    65
                        79 100
                                 90 100
                                                       93 100
                                                                64
                                                                         98
                                                                              45
                                                                                  66
                                                                                      55
[20]
      34
               59 100
                        73
                            21
                                 56 100
                                          36
                                              45
                                                  62
                                                       45
                                                           26
                                                                47
                                                                         94
                                                                             82 100 100
                   38
73
                                                               51
[39]
      89
           87 100
                        54
                             94
                                 55
                                     56
                                          74
                                              97
                                                   68 100
                                                            38
                                                                    41
                                                                         94 100
                                                                                  45
                                                                                      54
[58]
                        93
                                              37
                                                                         72
          52
72
                            58 100
                                      89
                                          53
      31
               60
                                                   29
                                                       23
                                                            89 100 100
                                                                             92
                                                                                  39
                                                                                      31
[77]
      80
```

The output gives the number of times each node was identified in bootstrapped analyses (the order is the same as in the original object). It is easily represented using nodelabels:

```
> plot(tre2, show.tip = FALSE, edge.width = 2)
> title("NJ tree + bootstrap values")
> tiplabels(frame = "none", pch = 20, col = transp(num2col(annot$year,
+ col.pal = myPal), 0.7), cex = 3, fg = "transparent")
> axisPhylo()
> temp <- pretty(1993:2008, 5)
> legend("topright", fill = transp(num2col(temp, col.pal = myPal),
+ 0.7), leg = temp, ncol = 2)
> nodelabels(myBoots, cex = 0.6)
```





As we can see, some nodes are very poorly supported. One common practice is to collapse these nodes into multifurcations. There is no dedicated method for this in *ape*, but one simple workaround consists in setting the corresponding edges to a length of zero (here, with bootstrap < 70%), and then collapsing the small branches:

```
> temp <- tre2
> N <- length(tre2$tip.label)
> toCollapse <- match(which(myBoots < 70) + N, temp$edge[, 2])
> temp$edge.length[toCollapse] <- 0
> tre3 <- di2multi(temp, tol = 1e-05)</pre>
```

The new tree might be slightly less informative, but more robust than the previous one:

```
> plot(tre3, show.tip = FALSE, edge.width = 2)
> title("NJ tree after collapsing weak nodes")
> tiplabels(annot$year, bg = transp(num2col(annot$year, col.pal = myPal),
+ 0.7), cex = 0.5, fg = "transparent")
> axisPhylo()
> temp <- pretty(1993:2008, 5)
> legend("topright", fill = transp(num2col(temp, col.pal = myPal),
+ 0.7), leg = temp, ncol = 2)
```





3 Maximum parsimony phylogenies

3.1 Introduction

Phylogenetic reconstruction based on parsimony seeks trees which minimize the total number of changes (substitutions) from ancestors to descendents. While a number of criticisms can be made to this approach, it is a simple way to infer phylogenies for data which display moderate to low divergence (i.e. most taxa differ from each other by only a few nucleotides, and the overall substitution rate is low).

In practice, there is often no way to perform an exhaustive search amongst all possible trees to find the most parsimonious one, and heuristic algorithms are used to browse the space of possible trees. The strategy is fairly simple: i) initialize the algorithm using a tree and ii) make small changes to the tree and retain those leading to better parsimony, until the parsimony score stops improving.

3.2 Implementation

Parsimony-based phylogenetic reconstruction is implemented in the package *phangorn*. It requires a tree (in *ape*'s format, i.e. a **phylo** object) and the original DNA sequences in *phangorn*'s own format, **phyDat**. We convert the data and generate a tree to initialize the method:

```
> dna2 <- as.phyDat(dna)
> class(dna2)
[1] "phyDat"
> dna2
80 sequences with 1701 character and 269 different site patterns.
The states are a c g t
> tre.ini <- nj(dist.dna(dna, model = "raw"))
> tre.ini
Phylogenetic tree with 80 tips and 78 internal nodes.
Tip labels:
        CY013200, CY013781, CY012128, CY013613, CY012160, CY012272, ...
Unrooted; includes branch lengths.
The parsimony of a given tree is given by:
> parsimony(tre.ini, dna2)
[1] 422
```

Then, optimization of the parsimony is achieved by:
> tre.pars <- optim.parsimony(tre.ini, dna2)
Final p-score 420 after 2 nni operations
> tre.pars
Phylogenetic tree with 80 tips and 78 internal nodes.
Tip labels:
 CY013200, CY013781, CY012128, CY013613, CY012160, CY012272, ...
Unrooted; no branch lengths.

Here, the final result is very close to the original tree. The obtained tree is unrooted and does not have branch lengths, but it can be plotted as previously:

```
> plot(tre.pars, type = "unr", show.tip = FALSE, edge.width = 2)
> title("Maximum-parsimony tree")
> tiplabels(annot$year, bg = transp(num2col(annot$year, col.pal = myPal),
+ 0.7), cex = 0.5, fg = "transparent")
> temp <- pretty(1993:2008, 5)
> legend("topright", fill = transp(num2col(temp, col.pal = myPal),
+ 0.7), leg = temp, ncol = 2, bg = transp("white"))
```



Maximum-parsimony tree

In this case, parsimony gives fairly consistent results with other approaches, which is only to be expected whenever the amount of divergence between the sequences is fairly low, as is the case in our data.

4 Maximum likelihood phylogenies

4.1 Introduction

Maximum likelihood phylogenetic reconstruction is somehow similar to parsimony methods in that it browses a space of possible tree topologies looking for the 'best' tree. However, it offers far more flexibility in that any model of sequence evolution can be taken into account. Given one model of evolution, one can compute the likelihood of a given tree, and therefore optimization procedures can be used to infer both the most likely tree topology and model parameters.

As in distance-based methods, model-based phylogenetic reconstruction requires thinking about which parameters should be included in a model. Usually, all possible substitutions are allowed to have different rates, and the substitution rate is allowed to vary across sites according to a gamma distribution. We refer to this model as $\text{GTR} + \Gamma(4)$ (GTR: global reversible time). More information about phylogenetic models can be found in [3].

4.2 Sorting out the data

Likelihood-based phylogenetic reconstruction is implemented in the package *phangorn*. Like parsimony-based approaches, it requires a tree (in *ape*'s format, i.e. a **phylo** object) and the original DNA sequences in *phangorn*'s own format, **phyDat**. As in the previous section, we convert the data and generate a tree to initialize the method:

Unrooted; includes branch lengths.

To initialize the optimization procedure, we need an initial fit for the model chosen. This is computed using pml:

```
> pml(tre.ini, dna2, k = 4)
```

loglikelihood: NaN

```
unconstrained loglikelihood: -4736.539
Discrete gamma model
Number of rate categories: 4
Shape parameter: 1
Rate matrix:
a c g t
a 0 1 1 1
c 1 0 1 1
g 1 1 0 1
t 1 1 1 0
Base frequencies:
0.25 0.25 0.25 0.25
```

The computed likelihood is NA, which is obviously a bit of a problem, but a likely frequent issue. This issue is due to missing data (NAs) in the original dataset. We therefore need to remove missing data before going further.

We first retrieve the position of missing data, i.e. any data differing from 'a', 'g', 'c' and 't'.

We can easily plot the number of missing data for each site:



The begining of the alignment is guilty for most of the missing data, which was only to be expected (extremities of the sequences have variable length).

```
> dna3 <- dna[, -na.posi]
> dna3
80 DNA sequences in binary format stored in a matrix.
All sequences of same length: 1596
Labels: CY013200 CY013781 CY012128 CY013613 CY012160 CY012272 ...
```

The object dna3 is an alignment of all sequences excluding missing data; dna4 is its conversion in phyDat format.

4.3 Getting a ML tree

We recompute the likelihood of the initial tree using pml:

```
> dna4 <- as.phyDat(dna3)
> tre.ini <- nj(dist.dna(dna3, model = "TN93"))
> fit.ini <- pml(tre.ini, dna4, k = 4)
> fit.ini
loglikelihood: -5183.648
unconstrained loglikelihood: -4043.367
Discrete gamma model
Number of rate categories: 4
Shape parameter: 1
Rate matrix:
    a c g t
    a 0 1 1 1
    c 1 0 1 1
    g 1 1 0 1
    t 1 1 1 0
Base frequencies:
    0.25 0.25 0.25 0.25
```

We now have all the information needed for seeking a maximum likelihood solution using optim.pml; we specify that we want to optimize tree topology (optNni=TRUE), base frequencies (optBf=TRUE), the rates of all possible subtitutions (optQ=TRUE), and use a gamma distribution to model variation in the substitution rates across sites (optGamma=TRUE):

```
> fit <- optim.pml(fit.ini, optNni = TRUE, optBf = TRUE, optQ = TRUE,
+ optGamma = TRUE)
```

> fit

```
loglikelihood: -4915.914
unconstrained loglikelihood: -4043.367
Discrete gamma model
Number of rate categories: 4
Shape parameter: 0.2843749
Rate matrix:
          a
              2.3831923 8.2953873
a 0.000000
                                       0.855505
              0.0000000 \ 0.1486215 \ 10.076469
c 2.383192
             0.1486215 0.0000000
 8.295387
                                      1.000000
g 8.295387 0.1400210 0.000000
t 0.855505 10.0764688 1.0000000
                                      0.000000
Base frequencies:
0.3416519 0.1953526 0.2242948 0.2387007
> class(fit)
[1] "pml"
> names(fit)
                             "k"
                                                     "Q"
                                                                "bf"
                 "inv"
 [1] "logLik"
                                         "shape"
                                                                            "rate"
 [8] "siteLik"
                 "weight"
                             "g"
                                         "w"
                                                     "eig"
                                                                "data"
                                                                            "model"
[15]
     "INV"
                 "11.0"
                             "tree"
                                         "lv"
                                                     "call"
                                                                "df"
                                                                            "wMix"
[22] "11Mix"
```

fit is a list with class pml storing various useful information about the model parameters and the optimal tree (stored in fit\$tree). In this example, we can see from the output that transitions $(a \leftrightarrow g \text{ and } c \leftrightarrow t)$ are much more frequent than transversions (other changes), which is consistent with biological expectations (transversions induce more drastic changes of chemical properties of the DNA and are more prone to purifying selection). We can verify that the optimized tree is indeed better than the original one using standard likelihood ration tests and AIC:

```
> anova(fit.ini, fit)
```

```
Likelihood Ratio Test Table

Log lik. Df Df change Diff log lik. Pr(>|Chi|)

1 -5183.6 158

2 -4915.9 166 8 535.47 < 2.2e-16

> AIC(fit.ini)

[1] 10683.3
```

> AIC(fit)

```
[1] 10163.83
```

Yes, the new tree is actually better than the initial one.

We can extract and plot the tree as we did before with other methods:

```
> tre4 <- root(fit$tree, 1)
> plot(tre4, show.tip = FALSE, edge.width = 2)
> title("Maximum-likelihood tree")
> tiplabels(annot$year, bg = transp(num2col(annot$year, col.pal = myPal),
+ 0.7), cex = 0.5, fg = "transparent")
> axisPhylo()
> temp <- pretty(1993:2008, 5)
> legend("topright", fill = transp(num2col(temp, col.pal = myPal),
+ 0.7), leg = temp, ncol = 2)
```



Maximum-likelihood tree

This tree is statistically better than the original NJ tree based on Tamura and Nei's distance [7]. However, we can note that it is remarkably similar to the 'robust' version of this distance-based tree (after collapsing weakly supported nodes). The structure of this dataset is fairly simple, and all methods give fairly consistent results. In practice, different methods can lead to different interpretations, and it is probably worth exploring different approaches before drawing conclusions on the data.

References

- [1] S. Dray and A.-B. Dufour. The ade4 package: implementing the duality diagram for ecologists. *Journal of Statistical Software*, 22(4):1–20, 2007.
- [2] T. Jombart. adegenet: a R package for the multivariate analysis of genetic markers. *Bioinformatics*, 24:1403–1405, 2008.
- [3] Scot A Kelchner and Michael A Thomas. Model use in phylogenetics: nine key questions. *Trends Ecol Evol*, 22(2):87–94, Feb 2007.
- [4] E. Paradis, J. Claude, and K. Strimmer. APE: analyses of phylogenetics and evolution in R language. *Bioinformatics*, 20:289–290, 2004.
- [5] R Development Core Team. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria, 2009. ISBN 3-900051-07-0.
- [6] Klaus Peter Schliep. phangorn: phylogenetic analysis in r. Bioinformatics, 27(4):592–593, Feb 2011.
- [7] K. Tamura and M. Nei. Estimation of the number of nucleotide substitutions in the control region of mitochondrial dna in humans and chimpanzees. *Mol Biol Evol*, 10(3):512–526, May 1993.