

732A51 Bioinformatics

LAB 2 Bioinformatics

Hugo Morvan
William Wiik

STIMA
Department of Computer and Information Science
Linköpings universitet

2024-11-25

Contents

1	Question 1	1
1.1	Question 1.1	2
1.2	Question 1.2 *	3
2	Question 2	4
2.1	Question 2.1	4
2.2	Question 2.2*	6
2.3	Question 2.3	6
3	Question 3	8
3.1	Question 3.1	8
3.2	Question 3.2*	12

1 Question 1

In this exercise you will perform statistical analysis of three nucleotide data sets. First download the sequences from GenBank and save them in a fasta file. For this use the provided R script, 732A51_BioinformaticsHT2023_Lab02_GenBankGetCode.R. This is a dataset of the RAG1 gene sequences from 33 lizard species. You are encouraged to read in detail the references in the script as they indicate many useful tools. Explore the dataset using the tools provided by the `ape` and `seqinr` packages. Take note of the lengths of all the sequences and the base composition.

```
## Gene bank accession numbers taken from http://www.jcsantosresearch.org/Class_2014_Spring_Comparative/
lizards_accession_numbers <- c("JF806202", "HM161150", "FJ356743", "JF806205",
                              "JQ073190", "GU457971", "FJ356741", "JF806207",
                              "JF806210", "AY662592", "AY662591", "FJ356748",
                              "JN112660", "AY662594", "JN112661", "HQ876437",
                              "HQ876434", "AY662590", "FJ356740", "JF806214",
                              "JQ073188", "FJ356749", "JQ073189", "JF806216",
                              "AY662598", "JN112653", "JF806204", "FJ356747",
                              "FJ356744", "HQ876440", "JN112651", "JF806215",
                              "JF806209")
lizards_sequences<-ape::read.GenBank(lizards_accession_numbers)
print(lizards_sequences)
```

```
## 33 DNA sequences in binary format stored in a list.
##
## Mean sequence length: 1982.879
##   Shortest sequence: 931
##   Longest sequence: 2920
##
## Labels:
## JF806202
## HM161150
## FJ356743
## JF806205
## JQ073190
## GU457971
## ...
##
## Base composition:
##   a   c   g   t
## 0.312 0.205 0.231 0.252
## (Total: 65.44 kb)
```

```
ape::write.dna(lizards_sequences, file = "lizard_seqs.fasta", format = "fasta",
               append = FALSE, nbcol = 6, colsep = " ", colw = 10)
```

1.1 Question 1.1

Simulate an artificial DNA sequence dataset. It should contain 33 sequence. The lengths of the sequences should be the same as in the lizard dataset, i.e. for each real sequence simulate an artificial one. The simulation rule is as follows, each nucleotide is to be independently and randomly drawn from the distribution given by the base composition (frequencies) in the true lizard sequences. Save your dataset in a fasta format file. Remember to give unique names to your sequences. Report on the base composition in your simulated data.

Answer:

```
true_base_composition <- c(0.312, 0.205, 0.231, 0.252) # a c g t

# Length of each real sequence
sequence_lengths <- sapply(lizards_sequences, length)

# Simulating DNA
sim_sequences <- lapply(sequence_lengths, function(length) {
  rDNABin(n = length, base.freq = true_base_composition)[[1]]
})

names(sim_sequences) <- paste0("sim_", names(lizards_sequences))

sim_sequences <- structure(sim_sequences, class = "DNABin")

ape::write.dna(sim_sequences, file = "sim_lizard_seqs.fasta",
               format = "fasta", append = FALSE, nbc = 6, colsep = " ", colw = 10)

sim_lizard_seqs <- read.FASTA("sim_lizard_seqs.fasta")
print(sim_lizard_seqs)
```

```
## 33 DNA sequences in binary format stored in a list.
##
## Mean sequence length: 1982.879
##   Shortest sequence: 931
##   Longest sequence: 2920
##
## Labels:
## sim_JF806202
## sim_HM161150
## sim_FJ356743
## sim_JF806205
## sim_JQ073190
## sim_GU457971
## ...
##
## Base composition:
##   a       c       g       t
## 0.312 0.203 0.230 0.255
## (Total: 65.44 kb)
```

a: 31.2%

c: 20.3%

g: 23.0%

t: 25.5%

We obtain almost the same base composition as the true base composition for lizard sequences, which is to be expected.

1.2 Question 1.2 *

2 Question 2

2.1 Question 2.1

Report some basic statistics on each sequence dataset: individual base composition, *GC* content, *CG*, *AT* content. Also translate your sequences into protein sequences (see Lab 1) and report on the amino acid composition. In your simulated sequences, how many times did you observe a stop codon inside your sequence? Does this occur in your true sequences? Comment.

```
true_sequences <- read.FASTA("lizard_seqs.fasta")
sim_sequences <- read.FASTA("sim_lizard_seqs.fasta")

calc_seq_stats <- function(all_seq, index_seq) {

  seq_DNABin <- structure(all_seq[[index_seq]], class = "DNABin")
  seq <- as.character(seq_DNABin) # from raw to character

  # How many g and c that are in the sequence
  GC_content <- ((sum(seq == "g") + sum(seq == "c")) / length(seq))

  # Base composition
  base <- base.freq(seq_DNABin) %>% round(3)

  # CpG (How many c is followed by a g)
  # seq[-length(seq)] removes the last element
  # seq[-1] removes the first element
  CpG <- sum(seq[-length(seq)] == "c" & seq[-1] == "g") / length(seq)

  # ApT (How many a is followed by a t)
  ApT <- sum(seq[-length(seq)] == "a" & seq[-1] == "t") / length(seq)

  print(paste0("Base composition for ", names(all_seq)[index_seq],
    ": a=", base["a"], ", c=", base["c"], ", g=", base["g"], ", t=", base["t"]))

  print(paste0("GC content for ", names(all_seq)[index_seq], ": ", round(GC_content, 3)))

  print(paste0("CpG content for ", names(all_seq)[index_seq], ": ", round(CpG, 3)))

  print(paste0("ApT content for ", names(all_seq)[index_seq], ": ", round(ApT, 3)))

}
```

```

# True lizard sequence 1 (JF806202)
calc_seq_stats(all_seq = true_sequences, index_seq = 1)

## [1] "Base composition for JF806202: a=0.29, c=0.203, g=0.244, t=0.264"
## [1] "GC content for JF806202: 0.446"
## [1] "CpG content for JF806202: 0.013"
## [1] "ApT content for JF806202: 0.077"

# Simulate lizard sequence 1 (sim_JF806202)
calc_seq_stats(all_seq = sim_sequences, index_seq = 1)

## [1] "Base composition for sim_JF806202: a=0.286, c=0.207, g=0.257, t=0.251"
## [1] "GC content for sim_JF806202: 0.464"
## [1] "CpG content for sim_JF806202: 0.045"
## [1] "ApT content for sim_JF806202: 0.07"

# "Cleanning" the sequences because they were some issues on my personal computer
dna_seqs <- readDNASTringSet("lizard_seqs.fasta")

## Warning in .Call2("fasta_index", filexp_list, nrec, skip, seek.first.rec, :
## reading FASTA file lizard_seqs.fasta: ignored 5451 invalid one-letter sequence
## codes

cleaned_seqs <- DNASTringSet(gsub("[^ATCG]", "", dna_seqs))

dna_seqs_sim <- readDNASTringSet("sim_lizard_seqs.fasta")

## Warning in .Call2("fasta_index", filexp_list, nrec, skip, seek.first.rec, :
## reading FASTA file sim_lizard_seqs.fasta: ignored 5451 invalid one-letter
## sequence codes

cleaned_seqs_sim <- DNASTringSet(gsub("[^ATCG]", "", dna_seqs_sim))

# File from Ying Luo
codon_and_ORF <- read.csv("codon_and_ORF.csv")

vec <- c()
for(i in 1:33){
  start_codon<- codon_and_ORF$codon_start[i]
  seq = cleaned_seqs[[i]]
  seq_vec = unlist(strsplit(tolower(as.character(seq)), ''))
  translation = seqinr::translate(seq_vec, start_codon-1)
  vec <- c(vec, translation)
}

char_table <- prop.table(table(vec))
round(char_table[order(table(vec))], 2)

```

```
## vec
##      W      *      Y      M      H      N      C      Q      F      G      I      P      D      T      A      R
## 0.01 0.01 0.02 0.03 0.03 0.03 0.03 0.03 0.04 0.05 0.05 0.05 0.05 0.05 0.06 0.06
##      V      E      S      K      L
## 0.06 0.07 0.08 0.08 0.10
```

```
vec_sim <- c()
for(i in 1:33){
  start_codon<- codon_and_ORF$codon_start[i]
  seq = cleaned_seqs_sim[[i]]
  seq_vec = unlist(strsplit(tolower(as.character(seq)), ''))
  translation = seqinr::translate(seq_vec)
  vec_sim <- c(vec_sim, translation)
}

char_table <- prop.table(table(vec_sim))
round(char_table[order(table(vec_sim))], 2)
```

```
## vec_sim
##      W      M      C      F      H      D      Q      Y      E      P      N      A      K      G      V      I
## 0.01 0.02 0.03 0.03 0.03 0.03 0.03 0.04 0.04 0.04 0.05 0.05 0.05 0.05 0.06 0.06
##      *      T      L      S      R
## 0.06 0.07 0.08 0.09 0.09
```

For the true sequences, we get almost no stopping codons (*), while for the simulated sequences, we get a proportion of stopping codons proportional to the number of ways to obtain a stopping codon from the standard codon table. This is to be expected since the simulated sequences generate proteins at random. For the true sequences, we expect to have much less stopping codons since they should only appear at the end of the protein.

2.2 Question 2.2*

2.3 Question 2.3

Align your sequences using software of your choice (a starter for R: <https://stackoverflow.com/questions/4497747/how-to-perform-basic-multiple-sequence-alignments-in-r>, you can also look what Biopython, BioPerl offer, use the Clustal family of programs or something else of your choice). Choose a distance measure between sequences, calculate for each alignment the distances between all pairs of sequences. Then, plot heatmaps visualizing the distances. Comment on what you can observe.

```
clust_lizard_seqs <- msaClustalW(cleaned_seqs, type = "dna")
```

```
## use default substitution matrix
```

```
clust_sim_lizard_seqs <- msaClustalW(cleaned_seqs_sim, type = "dna")
```

```
## use default substitution matrix
```

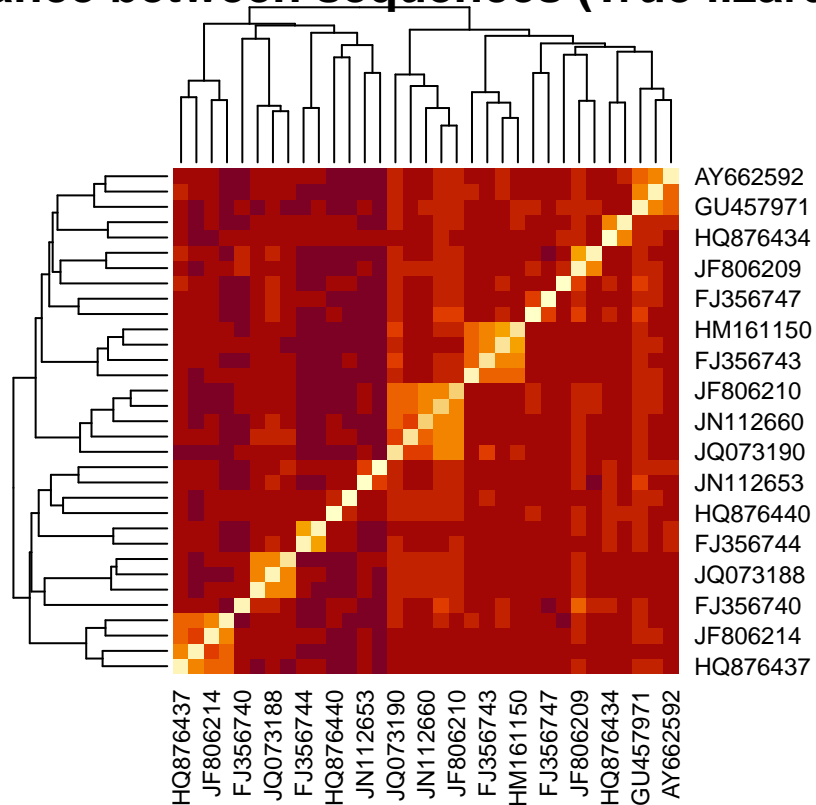


```
# Convert MSA object to formats used in other sequence
# type = "seqinr::alignment" because we want to use dist.alignment() later
converted_lizard_seqs <- msaConvert(clust_lizard_seqs, type = "seqinr::alignment")
converted_sim_lizard_seqs <- msaConvert(clust_sim_lizard_seqs, type = "seqinr::alignment")

# Calculate pairwise distances using dist.alignment() from seqinr package
dist_matrix_lizard_seqs <- as.matrix(dist.alignment(converted_lizard_seqs), matrix = "identity")
dist_matrix_sim_lizard_seqs <- as.matrix(dist.alignment(converted_sim_lizard_seqs), matrix = "identity")

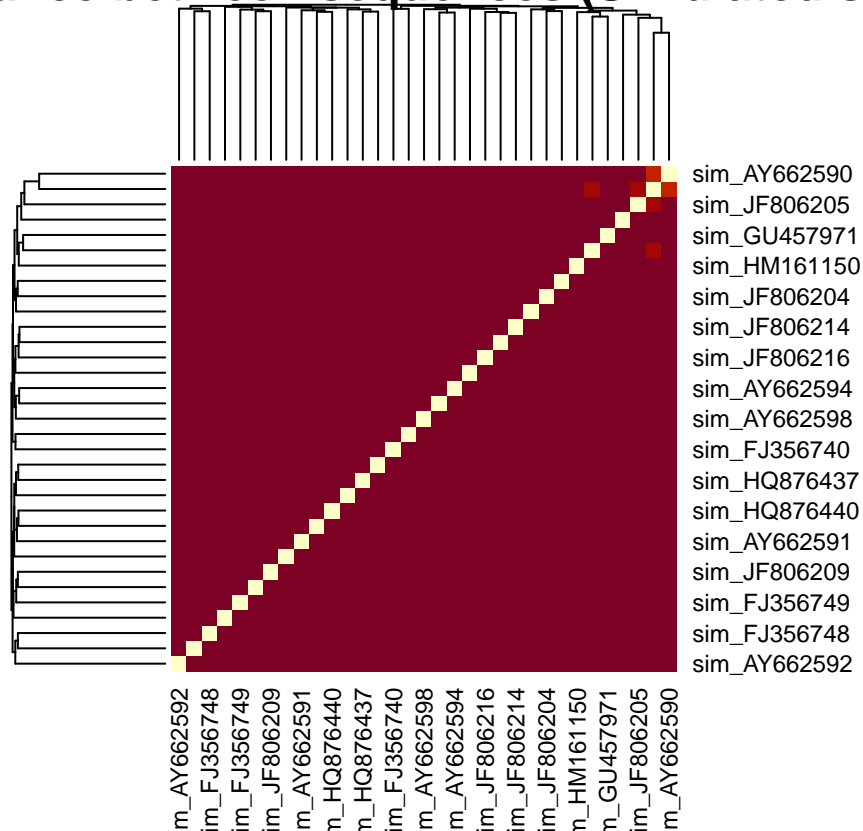
heatmap(dist_matrix_lizard_seqs, main = "Distance between sequences (True lizard seq.)")
```

Distance between sequences (True lizard seq.)



```
heatmap(dist_matrix_sim_lizard_seqs, main = "Distance between sequences (Simulated seq.)")
```

Distance between sequences (Simulated seq.)



Answer

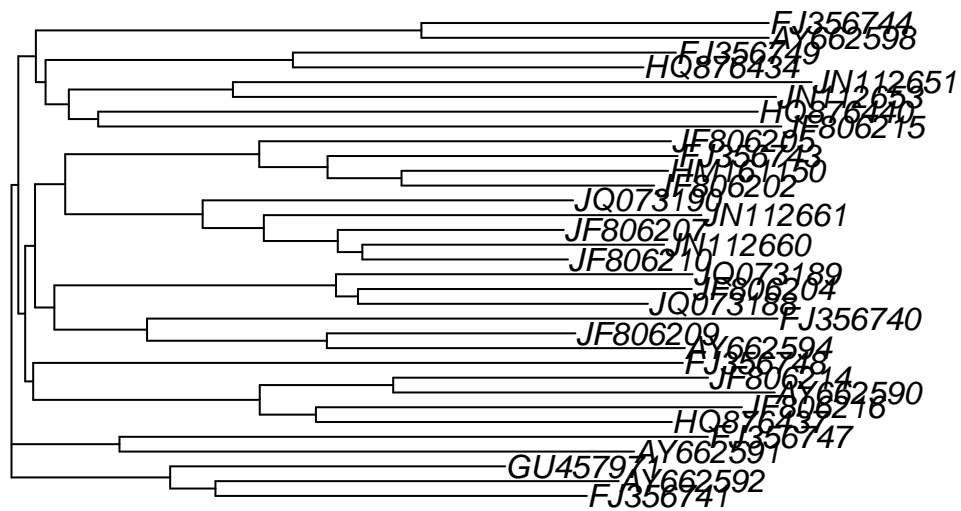
We observe some similarities between genes in the true data (heat on the heatmap), while for the simulated data, there is practically no similarity. This is expected because the simulated data comes from a iid process, while the true data comes from genetic evolution.

3 Question 3

3.1 Question 3.1

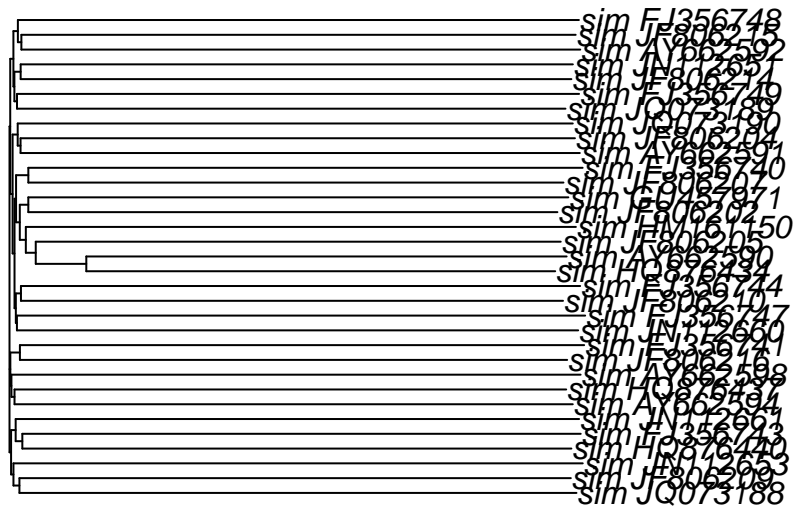
Construct (using algorithm and software of your choice) phylogenetic trees from the three multiple alignments (or distance matrices) done in Question 2.3. You might want to look at the functions offered by **ape**, **phangorn** (<https://cran.r-project.org/web/packages/phangorn/vignettes/Trees.pdf>) or go for some completely different software. Plot the inferred trees. Are the two based on the simulated data similar to expected? Perform a phylogenetic bootstrap analysis and report the bootstrap support for the individual clades, you can look at `ape::boot.phylo()`.

```
# Neighbor-Joining Tree Estimation :
phyl1 <- nj(dist_matrix_lizard_seqs)
plot(phyl1)
```



```
#Other options to create the tree : fastme.bal bionj, mvr(X, V), nj
```

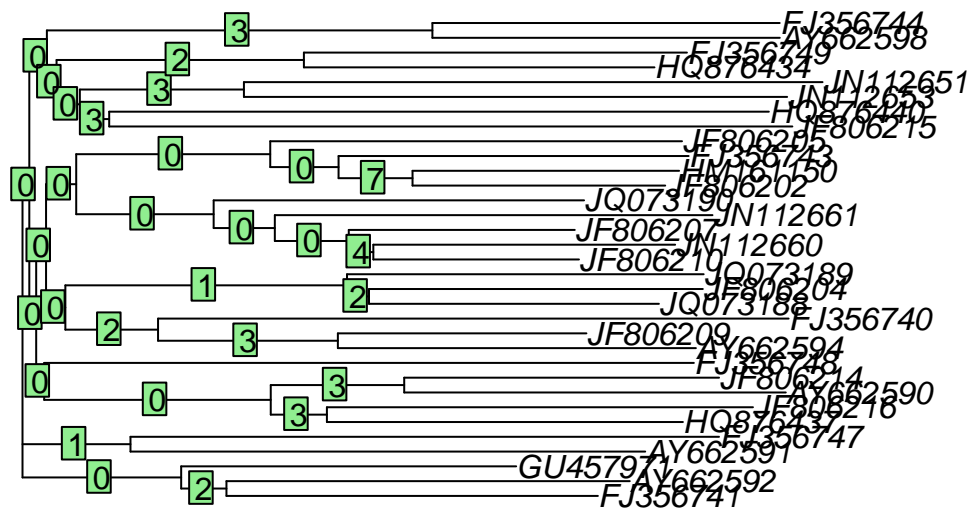
```
# Neighbor-Joining Tree Estimation :
phyl2 <-nj(dist_matrix_sim_lizard_seqs)
plot(phyl2)
```



```
#Boostrapping
## from ?boot.phylo:
f <- function(x) nj(x)
tw <- f(dist_matrix_lizard_seqs)
## bootstrap with 1000 replications:
(bp <- boot.phylo(tw, dist_matrix_lizard_seqs, f, B=100, quiet=TRUE))

## [1] NA 0 0 0 0 0 0 0 0 0 0 3 1 2 0 0 2 3 0 0 0 0 4 2 0 3
## [26] 3 1 2 7 3 3

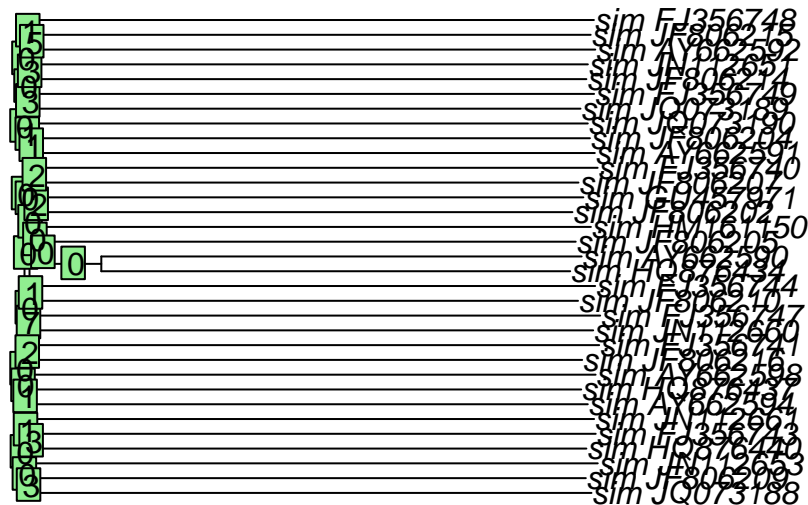
## more readable but the tree is really unrooted:
plot(tw)
drawSupportOnEdges(bp)
```



```
#Bootstrapping
## from ?boot.phylo:
f <- function(x) nj(x)
tw <- f(dist_matrix_sim_lizard_seqs)
## bootstrap with 100 replications:
(bp <- boot.phylo(tw, dist_matrix_sim_lizard_seqs, f, B=100, quiet=TRUE))

## [1] NA 0 0 0 0 0 0 0 0 0 0 0 0 7 1 3 1 0 1 1 1 1 3 3 5 2
## [26] 0 3 2 2 0 0

plot(tw)
drawSupportOnEdges(bp)
```



Interpretation Scale:

70-95%: Moderately supported branch 95-100%: Highly supported branch <70%: Weak support, less confidence in that branching pattern

For both true and simulated trees, none of the branches have significant support (>70%, i.e. support>70) Tree building functions tested : fastme.bal, bionj, nj., all yielded the same results.

3.2 Question 3.2*

Compare your inferred trees and also your simulated one. Apart from visualizing the trees one may calculate various indices related to them and distances between the trees. Explore what indices and metrics the **ape**, **distory**, **phangorn**, **phyloTop**, **TotalCopheneticIndex** or **treospace** R packages offer, choose some and report the results in a meaningful way. You might have to save your tree to drive and then read in it using e.g. **ape**'s tree reading functionality.

```
ape::comparePhylo(phy11, phy12)
```

```
## => Comparing phy11 with phy12.
```

```
## Both trees have the same number of tips: 33.
```

```
## Tips in phy11 not in phy12 : JF806210, JN112660, JN112661, JF806202, HM161150, FJ356743, JF806205, JQ
```

```
## Tips in phyl2 not in phyl1 : sim_HQ876434, sim_AY662590, sim_JF806205, sim_JF806202, sim_GU457971, si
## Both trees have the same number of nodes: 31.
## Both trees are unrooted.
## Both trees are not ultrametric.
```

```
#Metrics :
```

```
branching.times(phyl1)
```

```
##      34      35      36      37      38      39      40
## 0.15316931 0.15179162 0.15040138 0.14878227 0.14840381 0.14823928 0.14627806
##      41      42      43      44      45      46      47
## 0.14449828 0.14230621 0.14157215 0.13562936 0.13132725 0.12575191 0.11449990
##      48      49      50      51      52      53      54
## 0.12107332 0.11191276 0.10837965 0.10212980 0.10305594 0.10295326 0.08716938
##      55      56      57      58      59      60      61
## 0.08220290 0.09624948 0.08923640 0.09159566 0.08937464 0.08752518 0.08308446
##      62      63      64
## 0.07427249 0.07595380 0.07025427
```

```
branching.times(phyl2)
```

```
##      34      35      36      37      38      39      40      41
## 0.4275893 0.4270439 0.4269598 0.4261432 0.4263924 0.4252625 0.4246524 0.4239261
##      42      43      44      45      46      47      48      49
## 0.4246124 0.4238175 0.4228314 0.4226558 0.4220246 0.4238257 0.4219514 0.4214440
##      50      51      52      53      54      55      56      57
## 0.4198596 0.4212946 0.4229515 0.4189981 0.4192971 0.4192182 0.4201420 0.4186683
##      58      59      60      61      62      63      64
## 0.4197654 0.4153275 0.4179461 0.4134634 0.4134462 0.4078837 0.3699204
```

```
# #compar.cheverud      Cheverud's Comparative Method
# #compar.gee          Comparative Analysis with GEEs
# #compar.lynch        Lynch's Comparative Method
# #compar.ou           Ornstein-Uhlenbeck Model for Continuous Characters
#
# a = compute.brlen(phyl1)      #Branch Lengths Computation
# plot(a)
# b = compute.brtime(phyl1)     #Compute and Set Branching Times
# plot(b)
#
# a = compute.brlen(phyl2)      #Branch Lengths Computation
# plot(a)
# b = compute.brtime(phyl2)     #Compute and Set Branching Times
# plot(b)
#
# heatmap(cophenetic.phylo(phyl1))
# heatmap(cophenetic.phylo(phyl2))
```

```

# corBlomberg           #Blomberg et al.'s Correlation Structure
# corBrownian           #Brownian Correlation Structure
# corClasses            #Phylogenetic Correlation Structures
# corGrafen             #Grafen's (1989) Correlation Structure
# corMartins            #Martins's (1997) Correlation Structure
# corPagel              #Pagel's "lambda" Correlation Structure
# corphylo              #Correlations among Multiple Traits with Phylogenetic Signal
# correlogram.formula   #Phylogenetic Correlogram
#
# dist.topo(phy1, phy2) #needs similar leaf names
#
# gammaStat             #Gamma-Statistic of Pybus and Harvey, requires ultrametric trees
#
# root                  #Roots Phylogenetic Trees

summary.phylo(phy1)

```

```

##
## Phylogenetic tree: phy1
##
##   Number of tips: 33
##   Number of nodes: 31
##   Branch lengths:
##     mean: 0.05191795
##     variance: 0.001511728
##     distribution summary:
##       Min.      1st Qu.      Median      3rd Qu.      Max.
## 0.001377696 0.014389981 0.051113051 0.075027992 0.138188147
##   No root edge.
##   First ten tip labels: JF806210
##                        JN112660
##                        JN112661
##                        JF806202
##                        HM161150
##                        FJ356743
##                        JF806205
##                        JQ073190
##                        JF806207
##                        JQ073188
##   No node labels.

```

```
summary.phylo(phy2)
```

```

##
## Phylogenetic tree: phy2
##
##   Number of tips: 33

```



```

## Number of nodes: 31
## Branch lengths:
##   mean: 0.2178854
##   variance: 0.04231989
##   distribution summary:
##           Min.      1st Qu.      Median      3rd Qu.      Max.
## 0.0005453872 0.0027112159 0.3607140035 0.4162849315 0.4275080112
## No root edge.
## First ten tip labels: sim_HQ876434
##                        sim_AY662590
##                        sim_JF806205
##                        sim_JF806202
##                        sim_GU457971
##                        sim_HM161150
##                        sim_JQ073189
##                        sim_JQ073190
##                        sim_JF806207
##                        sim_JN112660
## No node labels.

```

Answer:

True data phylogenetic tree Branch lengths stats:

mean: 0.05191795 , variance: 0.001511728

Min.: 0.001377696

1st Qu.: 0.014389981

Median: 0.051113051

3rd Qu.: 0.075027992

Max.: 0.138188147

Simulated data phylogenetic tree Branch lengths stats:

mean: 0.2178854 , variance: 0.04231989

Min.: 0.0005453872

1st Qu.: 0.0027112159

Median: 0.3607140035

3rd Qu. : 0.4162849315

Max.: 0.4275080112

The simulated data phylogenetic tree branch length has on average much longer branch lengths (0.05 vs 0.22), but with also much more variance (0.0015 vs 0.0423). These results correlate with the observation from the heatmaps which shows that in the simulated data, no clade is close to any other, which makes sense since the data was randomly generated and should be i.i.d.