

Bioinformatics - Computer Lab 2

Group 7: Lennart Schilling (*lensc874*), Thijs Quast (*thiqu264*), Mariano Maquieira Mariani (*marma330*)

27 November 2018

Question 1

At first, the dataset of the RAG1 gene sequences from 33 lizard species were downloaded from GenBank and saved in a fasta file using the provided R script *732A51 BioinformaticsHT2018 Lab02 GenBankGetCode.R*. The code can be found in Appendix 1 (Data Import of original dataset).

Question 1.1

The saved fasta-file has to be read in R so that we can work with that. The R code for the reading process can be found in Appendix 1.1 (Reading original data).

After that, the artificial dataset is built by considering that it contains 33 sequences (each length of the sequences is the same as in the lizard dataset) so that for each real sequence an artificial one is created. As mentioned, the simulation of the artificial sequences is based on the distribution given by the base composition of the original dataset.

The artificial dataset is submitted as the fasta file *artificial_dataset_1_1.fasta*. The written function for all these processes automatically prints the base composition in the simulated data compared to the base composition in the original data. An extract from the output can be seen here:

```
get_artificial_sequence_dataset(lizards_sequences)

## [1] "comparison of base compositions between original and artificial datasets (values rounded):"
##   name_original name_artificial a_original a_artificial c_original
## 1 "JF806202"      "1"           "0.29"      "0.29"      "0.2"
## 2 "HM161150"      "2"           "0.31"      "0.32"      "0.21"
## 3 "FJ356743"      "3"           "0.31"      "0.3"       "0.21"
## 4 "JF806205"      "4"           "0.28"      "0.29"      "0.21"
## 5 "JQ073190"      "5"           "0.31"      "0.29"      "0.2"
##   c_artificial g_original g_artificial t_original t_artificial
## 1 "0.2"         "0.24"      "0.24"      "0.26"      "0.27"
## 2 "0.22"        "0.23"      "0.22"      "0.24"      "0.24"
## 3 "0.21"        "0.23"      "0.24"      "0.24"      "0.25"
## 4 "0.21"        "0.24"      "0.25"      "0.26"      "0.26"
## 5 "0.21"        "0.24"      "0.23"      "0.26"      "0.27"
```

It becomes clear that the base compositions are very similar. The entire code for the function can be seen in Appendix 1.1 (Function code).

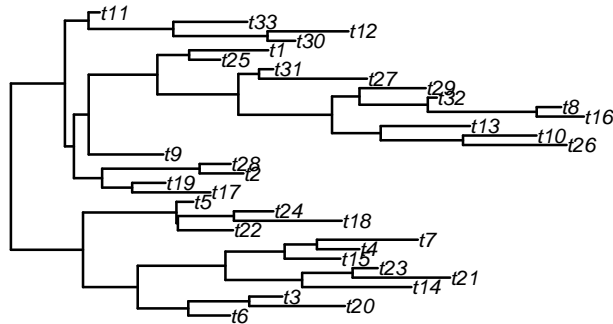
Question 1.2

In this part of the exercise do we use the prepared data from part 1, in Appendix 1 code can be found in (Data Import of original dataset).

We used the function *rtree* to create a tree object of the type phylo and the length of the original sequences.

```
tree <- rtree(n = length(lizards_sequences))
```

Here you can find the plot of the tree.



After the simulation of the phylogenetic tree, we had to simulate the sequence.

For this, we had several things to do. 1. We simulated a transition rate matrix (Q-Matrix). In this case we choose one by yourself.

2. We had to choose the lengths of the sequences. We chose the average length of the original sequences and used this length for every artificial sequence.

```
# calculating average length of original sequences
avg_length = c()
for (seq in 1:33) {
  avg_length = c(avg_length, length(lizards_sequences[[seq]]))
}
avg_length = mean(avg_length)
```

Now we can simulate the sequences by using the function *phangorn::simSeq()*.

```
sequences_artificial <- simSeq(tree, l = avg_length, Q=transition_matrix , type = "DNA")
```

Since in sequences are filled with integers from 1 to 4, do we have to replace the numbers by the letters a,b,c,d.

1 = a

2 = b

3 = c

4 = d

The code for this can be found in Appendix 1.2

The second simulate a artificial DNA sequence dataset do we save as “artificial_dataset_1_2.fasta”.

```
ape::write.dna(sequences_artificial, file = "artificial_dataset_1_2.fasta", format = "fasta", colsep = "
```

Question 2

Question 2.1

```
lizards_sequences = read.fasta("lizard_seqs.fasta")
original_dataset <- lizards_sequences
artificial_sequences_1 <- read.fasta("artificial_dataset_1_1.fasta")
artificial_sequences_2 <- read.fasta("artificial_dataset_1_2.fasta")
original_base_compositions <- list()
artificial_1_base_compositions <- list()
artificial_2_base_compositions <- list()
for (i in 1:length(original_dataset)) {
  # getting base compositions for each original sequence
  original_base_compositions[[i]] =
    seqinr::count(original_dataset[[i]],1)
}

for (i in 1:length(artificial_sequences_1)) {
  # getting base compositions for each original sequence
  artificial_1_base_compositions[[i]] =
    seqinr::count(artificial_sequences_1[[i]],1)
}

for (i in 1:length(artificial_sequences_2)) {
  # getting base compositions for each original sequence
  artificial_2_base_compositions[[i]] =
    seqinr::count(artificial_sequences_2[[i]],1)
}

Reduce('+', original_base_compositions)

##
##      a      c      g      t
## 20414 13422 15089 16474
sum(Reduce('+', original_base_compositions))

## [1] 65399
Reduce('+', original_base_compositions)/sum(Reduce('+', original_base_compositions))

##
##      a      c      g      t
## 0.3121454 0.2052325 0.2307222 0.2518999
Reduce('+', artificial_1_base_compositions)

##
##      a      c      g      t
## 20491 13346 15157 16441
sum(Reduce('+', artificial_1_base_compositions))

## [1] 65435
Reduce('+', artificial_1_base_compositions)/sum(Reduce('+', artificial_1_base_compositions))
```

```
##
##      a      c      g      t
## 0.3131505 0.2039581 0.2316344 0.2512570
Reduce('+', artificial_2_base_compositions)

##
##      a      c      g      t
## 16467 16604 16232 16103
sum(Reduce('+', artificial_2_base_compositions))

## [1] 65406
Reduce('+', artificial_2_base_compositions)/sum(Reduce('+', artificial_2_base_compositions))

##
##      a      c      g      t
## 0.2517659 0.2538605 0.2481730 0.2462007
```

The original dataset and the first artificially created dataset are rather similar in their distributions for A, C, T and G's. However, the second artificially created dataset has a slightly different distribution. This final dataset has almost uniform distribution for A, C, T and G's, they all occur with an average frequency of approximately 25%.

```
library(rDNase)
original_compositions <- list()
for (i in 1:length(lizards_sequences)) {
  string1 <- paste(lizards_sequences[[i]], collapse = "")
  string1 <- toupper(string1)
  original_compositions[[i]] <- kmer(string1)
}

artificial_compositions_1 <- list()
for (i in 1:length(artificial_sequences_1)) {
  string1 <- paste(artificial_sequences_1[[i]], collapse = "")
  string1 <- toupper(string1)
  artificial_compositions_1[[i]] <- kmer(string1)
}

artificial_compositions_2 <- list()
for (i in 1:length(artificial_sequences_2)) {
  string1 <- paste(artificial_sequences_2[[i]], collapse = "")
  string1 <- toupper(string1)
  artificial_compositions_2[[i]] <- kmer(string1)
}
```

```
Reduce('+', original_compositions)
```

```
##      TG      GA      AA      AG      AT      CA      TT      CT      TC
## 6960    4554    5664    4867    5096    3909    2277    4296    4947
##      CC      GC      GT      AC      GG (Other)  NA's
## 3207    3076    3172    2727    3237    4728    2685
```

```
Reduce('+', artificial_compositions_1)
```

```
##      AA      AC      AG      AT      CA      CC      CG      CT      GA      GC      GG      GT      TA      TC      TG
## 6407 4188 4677 5207 4137 2728 3116 3361 4680 3209 3520 3743 5252 3216 3839
##      TT
```

```
## 4122
```

```
Reduce('+', artificial_compositions_2)
```

```
##   AA   AC   AG   AT   CA   CC   CG   CT   GA   GC   GG   GT   TA   TC   TG
## 4162 4165 4099 4029 4149 4250 4068 4131 4099 4082 4036 4006 4048 4096 4022
##   TT
## 3931
```

GC content is the largest for the second artificially created dataset. CG content is largest for the second artificially created dataset. AT content is largest in the original dataset.

```
# Protein sequences
```

```
protein_original <- read.fasta("lizard_protein.fasta")
protein_artificial_1 <- read.fasta("artificial_1_protein.fasta")
protein_artificial_2 <- read.fasta("artificial_2_protein.fasta")
```

```
library(protr)
original_aac <- list()
for (i in 1:length(protein_original)) {
  string1 <- paste(protein_original[[i]], collapse = "")
  string1 <- toupper(string1)
  string1 <- gsub(pattern = "[*]", replacement = "", x = string1)
  string1 <- gsub(pattern = "B", replacement = "", x = string1)
  string1 <- gsub(pattern = "J", replacement = "", x = string1)
  string1 <- gsub(pattern = "O", replacement = "", x = string1)
  string1 <- gsub(pattern = "U", replacement = "", x = string1)
  string1 <- gsub(pattern = "X", replacement = "", x = string1)
  string1 <- gsub(pattern = "Z", replacement = "", x = string1)
  original_aac[[i]] <- extractAAC(string1)
}
```

```
artificial_1_aac <- list()
for (i in 1:length(protein_artificial_1)) {
  string1 <- paste(protein_artificial_1[[i]], collapse = "")
  string1 <- toupper(string1)
  string1 <- gsub(pattern = "[*]", replacement = "", x = string1)
  string1 <- gsub(pattern = "B", replacement = "", x = string1)
  string1 <- gsub(pattern = "J", replacement = "", x = string1)
  string1 <- gsub(pattern = "O", replacement = "", x = string1)
  string1 <- gsub(pattern = "U", replacement = "", x = string1)
  string1 <- gsub(pattern = "X", replacement = "", x = string1)
  string1 <- gsub(pattern = "Z", replacement = "", x = string1)
  artificial_1_aac[[i]] <- extractAAC(string1)
}
```

```
artificial_2_aac <- list()
for (i in 1:length(protein_artificial_2)) {
  string1 <- paste(protein_artificial_2[[i]], collapse = "")
  string1 <- toupper(string1)
  string1 <- gsub(pattern = "[*]", replacement = "", x = string1)
  string1 <- gsub(pattern = "B", replacement = "", x = string1)
  string1 <- gsub(pattern = "J", replacement = "", x = string1)
  string1 <- gsub(pattern = "O", replacement = "", x = string1)
  string1 <- gsub(pattern = "U", replacement = "", x = string1)
  string1 <- gsub(pattern = "X", replacement = "", x = string1)
  string1 <- gsub(pattern = "Z", replacement = "", x = string1)
}
```

```
artificial_2_aac[[i]] <- extractAAC(string1)
}
```

```
Reduce('+', original_aac)/length(original_aac)
```

```
##           A           R           N           D           C           E
## 0.04504567 0.06892454 0.03459927 0.03899492 0.04615282 0.06251367
##           Q           G           H           I           L           K
## 0.04665052 0.05212339 0.03881790 0.03970175 0.09512420 0.06888196
##           M           F           P           S           T           W
## 0.02177691 0.04072228 0.06041408 0.09363222 0.05396491 0.02155234
##           Y           V
## 0.02414733 0.04625931
```

```
Reduce('+', artificial_1_aac)/length(artificial_1_aac)
```

```
##           A           R           N           D           C           E
## 0.05057534 0.08796184 0.04476749 0.03495967 0.03066037 0.04049515
##           Q           G           H           I           L           K
## 0.04055862 0.05711907 0.02888776 0.06421406 0.09216594 0.05237905
##           M           F           P           S           T           W
## 0.01695355 0.03323242 0.04400956 0.09865235 0.06662918 0.01463302
##           Y           V
## 0.03807384 0.06307173
```

```
Reduce('+', artificial_2_aac)/length(artificial_2_aac)
```

```
##           A           R           N           D           C           E
## 0.06681774 0.09424995 0.03355575 0.03320151 0.03318094 0.03266863
##           Q           G           H           I           L           K
## 0.03317184 0.06386547 0.03403190 0.04785987 0.09575508 0.03288863
##           M           F           P           S           T           W
## 0.01713940 0.03254250 0.06717763 0.10010313 0.06718395 0.01835640
##           Y           V
## 0.03379739 0.06245228
```

After removing some unwanted letters and characters, the observed amino acids remain for the obtained protein sequences. Distribution of the amino acids among the three databases of obtained protein sequences is rather similar for all three protein databases.

```
library(seqinr)
library(stringr)

# reading original_dataset from fasta file
lizards_sequences = read.fasta("lizard_seqs.fasta")
# preparing data in fasta file (dna sequences include empty spaces which will be removed)
for (i in 1:length(lizards_sequences)) {
  lizards_sequences[[i]] = lizards_sequences[[i]][lizards_sequences[[i]] != " "]
}
taa_count <- c()
tag_count <- c()
tga_count <- c()

for (i in 1:33){
  string <- lizards_sequences[[i]]
  string <- paste(lizards_sequences[[i]], collapse = "")
```

```

taa_count[i] <- str_count(string, pattern = "taa")
tag_count[i] <- str_count(string, pattern = "tag")
tga_count[i] <- str_count(string, pattern = "tga")
}

names_sequences <- names(lizards_sequences)
df_original <- as.data.frame(cbind(names_sequences, taa_count, tag_count, tga_count,
                                   total_count_1 = taa_count + tag_count + tga_count))

artificial_sequences_1 <- read.fasta("artificial_dataset_1_1.fasta")
taa_a1 <- c()
tag_a1 <- c()
tga_a1 <- c()
for (i in 1:33){
  string <- artificial_sequences_1[[i]]
  string <- paste(artificial_sequences_1[[i]], collapse = "")
  taa_a1[i] <- str_count(string, pattern = "taa")
  tag_a1[i] <- str_count(string, pattern = "tag")
  tga_a1[i] <- str_count(string, pattern = "tga")
}

names_a1 <- names(artificial_sequences_1)

df_a1 <- as.data.frame(cbind(names_a1, taa_a1, tag_a1, tga_a1, total_count_2 =
                              taa_a1 + tag_a1 + tga_a1))

artificial_sequences_2 <- read.fasta("artificial_dataset_1_2.fasta")
taa_a2 <- c()
tag_a2 <- c()
tga_a2 <- c()

for (i in 1:33){
  string <- artificial_sequences_2[[i]]
  string <- paste(artificial_sequences_2[[i]], collapse = "")
  taa_a2[i] <- str_count(string, pattern = "taa")
  tag_a2[i] <- str_count(string, pattern = "tag")
  tga_a2[i] <- str_count(string, pattern = "tga")
}

names_a2 <- names(artificial_sequences_2)

df_a2 <- as.data.frame(cbind(names_a2, taa_a2, tag_a2, tga_a2, total_count_3 =
                              taa_a2 + tag_a2 + tga_a2))

df_all <- as.data.frame(cbind(df_a1, df_a2))
df_all

```

##	names_a1	taa_a1	tag_a1	tga_a1	total_count_2	names_a2	taa_a2	tag_a2
## 1	1	23	15	25	63	1	29	20
## 2	2	74	46	33	153	2	32	25
## 3	3	64	48	44	156	3	35	30
## 4	4	16	17	24	57	4	34	32
## 5	5	32	27	37	96	5	28	26
## 6	6	34	15	23	72	6	25	34

## 7	7	85	48	45	178	7	29	30
## 8	8	22	21	19	62	8	18	33
## 9	9	18	26	16	60	9	37	27
## 10	10	77	54	51	182	10	32	35
## 11	11	81	45	53	179	11	26	27
## 12	12	74	54	57	185	12	39	35
## 13	13	71	54	47	172	13	39	23
## 14	14	70	39	49	158	14	40	30
## 15	15	78	48	60	186	15	22	30
## 16	16	29	18	16	63	16	24	23
## 17	17	26	18	15	59	17	28	35
## 18	18	65	43	55	163	18	31	32
## 19	19	54	57	57	168	19	30	26
## 20	20	24	19	22	65	20	25	34
## 21	21	44	26	30	100	21	27	34
## 22	22	67	53	47	167	22	20	22
## 23	23	41	25	26	92	23	33	27
## 24	24	26	15	22	63	24	31	27
## 25	25	87	45	57	189	25	39	26
## 26	26	64	48	39	151	26	38	34
## 27	27	26	15	22	63	27	24	46
## 28	28	78	51	54	183	28	39	17
## 29	29	70	68	50	188	29	34	30
## 30	30	21	13	20	54	30	35	30
## 31	31	65	44	56	165	31	39	26
## 32	32	22	22	14	58	32	27	34
## 33	33	29	17	11	57	33	27	35
##	tga_a2	total_count_3						
## 1	29		78					
## 2	25		82					
## 3	28		93					
## 4	37		103					
## 5	27		81					
## 6	28		87					
## 7	32		91					
## 8	24		75					
## 9	29		93					
## 10	35		102					
## 11	36		89					
## 12	31		105					
## 13	28		90					
## 14	30		100					
## 15	29		81					
## 16	26		73					
## 17	30		93					
## 18	30		93					
## 19	38		94					
## 20	22		81					
## 21	24		85					
## 22	33		75					
## 23	31		91					
## 24	29		87					
## 25	34		99					
## 26	28		100					


```
## 27      28      98
## 28      32      88
## 29      34      98
## 30      36     101
## 31      35     100
## 32      26      87
## 33      23      85
```

Interpreting stop codons as either “taa”, “tag” or “tga” results in many stop codons for each sequence. In the original dataset this is highly unlikely, as a natural translation starts at a start codon and then continues until it reaches a stop codon. Or if it does not reach a stop codon at all.

Question 2.2

```
library(markovchain)
mcFitMle_original <- markovchainFit(lizards_sequences, method = "mle")
mcFitMle_original

## $estimate
## MLE Fit
## A 8 - dimensional discrete Markov Chain defined by the following states:
## a, c, g, m, r, s, t, y
## The transition matrix (by rows) is defined as follows:
##      a      c      g      m      r      s
## a 0.3377604 0.1730948 0.27493261 4.900760e-05 0.0002450380 0.000000e+00
## c 0.3793901 0.2477071 0.05010812 0.000000e+00 0.0003728283 0.000000e+00
## g 0.3934372 0.2029168 0.19323832 6.629102e-05 0.0003314551 0.000000e+00
## m 0.0000000 0.0000000 0.66666667 0.000000e+00 0.0000000000 0.000000e+00
## r 0.4117647 0.1764706 0.11764706 0.000000e+00 0.0000000000 0.000000e+00
## s 0.0000000 1.0000000 0.00000000 0.000000e+00 0.0000000000 0.000000e+00
## t 0.1508047 0.2115396 0.35718190 6.073489e-05 0.0001214698 6.073489e-05
## y 0.3333333 0.2000000 0.13333333 0.000000e+00 0.0000000000 0.000000e+00
##      t      y
## a 0.2136731 0.0002450380
## c 0.3222728 0.0001491313
## g 0.2096122 0.0003977461
## m 0.3333333 0.0000000000
## r 0.2941176 0.0000000000
## s 0.0000000 0.0000000000
## t 0.2801093 0.0001214698
## y 0.3333333 0.0000000000
##
##
## $standardError
##      a      c      g      m      r
## a 0.004068516 0.002912552 0.003670666 4.900760e-05 1.095843e-04
## c 0.005318784 0.004297725 0.001932963 0.000000e+00 1.667339e-04
## g 0.005106990 0.003667637 0.003579101 6.629102e-05 1.482312e-04
## m 0.000000000 0.000000000 0.471404521 0.000000e+00 0.000000e+00
## r 0.155632430 0.101885342 0.083189033 0.000000e+00 0.000000e+00
## s 0.000000000 1.000000000 0.000000000 0.000000e+00 0.000000e+00
## t 0.003026402 0.003584388 0.004657618 6.073489e-05 8.589211e-05
## y 0.149071198 0.115470054 0.094280904 0.000000e+00 0.000000e+00
##      s      t      y
```

```

## a 0.000000e+00 0.003235986 1.095843e-04
## c 0.000000e+00 0.004902089 1.054518e-04
## g 0.000000e+00 0.003727654 1.623792e-04
## m 0.000000e+00 0.333333333 0.000000e+00
## r 0.000000e+00 0.131533410 0.000000e+00
## s 0.000000e+00 0.000000000 0.000000e+00
## t 6.073489e-05 0.004124610 8.589211e-05
## y 0.000000e+00 0.149071198 0.000000e+00
##
## $confidenceLevel
## [1] 0.95
##
## $lowerEndpointMatrix
##      a      c      g m      r s      t
## a 0.33106824 0.168304107 0.26889491 0 6.478782e-05 0 0.20835040
## c 0.37064143 0.240637977 0.04692868 0 9.857546e-05 0 0.31420954
## g 0.38503694 0.196884079 0.18735122 0 8.763643e-05 0 0.20348075
## m 0.00000000 0.000000000 0.00000000 0 0.000000e+00 0 0.00000000
## r 0.15577214 0.008884115 0.00000000 0 0.000000e+00 0 0.07776444
## s 0.00000000 0.000000000 0.00000000 0 0.000000e+00 0 0.00000000
## t 0.14582675 0.205643836 0.34952080 0 0.000000e+00 0 0.27332494
## y 0.08813303 0.010068663 0.00000000 0 0.000000e+00 0 0.08813303
##      y
## a 6.478782e-05
## c 0.000000e+00
## g 1.306561e-04
## m 0.000000e+00
## r 0.000000e+00
## s 0.000000e+00
## t 0.000000e+00
## y 0.000000e+00
##
## $upperEndpointMatrix
##      a      c      g      m      r      s
## a 0.3444525 0.1778856 0.28097032 0.0001296179 0.0004252881 0.0000000000
## c 0.3881387 0.2547762 0.05328756 0.0000000000 0.0006470811 0.0000000000
## g 0.4018374 0.2089495 0.19912541 0.0001753300 0.0005752738 0.0000000000
## m 0.0000000 0.0000000 1.00000000 0.0000000000 0.0000000000 0.0000000000
## r 0.6677573 0.3440571 0.25448084 0.0000000000 0.0000000000 0.0000000000
## s 0.0000000 1.0000000 0.00000000 0.0000000000 0.0000000000 0.0000000000
## t 0.1557827 0.2174354 0.36484300 0.0001606349 0.0002627497 0.0001606349
## y 0.5785336 0.3899313 0.28841162 0.0000000000 0.0000000000 0.0000000000
##      t      y
## a 0.2189958 0.0004252881
## c 0.3303360 0.0003225840
## g 0.2157436 0.0006648361
## m 0.8816179 0.0000000000
## r 0.5104709 0.0000000000
## s 0.0000000 0.0000000000
## t 0.2868937 0.0002627497
## y 0.5785336 0.0000000000
mcFitMle_a1 <- markovchainFit(artificial_sequences_1, method = "mle")
mcFitMle_a1

```

```

## $estimate
## MLE Fit
## A 4 - dimensional discrete Markov Chain defined by the following states:
## a, c, g, t
## The transition matrix (by rows) is defined as follows:
##      a      c      g      t
## a 0.3128571 0.2045022 0.2283803 0.2542605
## c 0.3100735 0.2044671 0.2335482 0.2519113
## g 0.3088701 0.2117872 0.2323126 0.2470301
## t 0.3196786 0.1957514 0.2336722 0.2508978
##
##
## $standardError
##      a      c      g      t
## a 0.003908576 0.003160055 0.003339450 0.003523587
## c 0.004820830 0.003914725 0.004183866 0.004345236
## g 0.004514950 0.003738651 0.003915628 0.004037755
## t 0.004411144 0.003451810 0.003771359 0.003907895
##
## $confidenceLevel
## [1] 0.95
##
## $lowerEndpointMatrix
##      a      c      g      t
## a 0.3064280 0.1993043 0.2228874 0.2484647
## c 0.3021439 0.1980279 0.2266663 0.2447640
## g 0.3014437 0.2056377 0.2258719 0.2403886
## t 0.3124229 0.1900737 0.2274688 0.2444699
##
## $upperEndpointMatrix
##      a      c      g      t
## a 0.3192861 0.2097000 0.2338732 0.2600562
## c 0.3180030 0.2109062 0.2404300 0.2590585
## g 0.3162965 0.2179368 0.2387532 0.2536716
## t 0.3269343 0.2014291 0.2398755 0.2573257
mcFitMle_a2 <- markovchainFit(artificial_sequences_2, method = "mle")
mcFitMle_a2

```

```

## $estimate
## MLE Fit
## A 4 - dimensional discrete Markov Chain defined by the following states:
## a, c, g, t
## The transition matrix (by rows) is defined as follows:
##      a      c      g      t
## a 0.2529322 0.2531146 0.2491036 0.2448496
## c 0.2499699 0.2560549 0.2450898 0.2488854
## g 0.2526660 0.2516181 0.2487826 0.2469334
## t 0.2514754 0.2544574 0.2498602 0.2442070
##
##
## $standardError
##      a      c      g      t
## a 0.003920606 0.003922018 0.003890819 0.003857454
## c 0.003880753 0.003927704 0.003842684 0.003872325

```

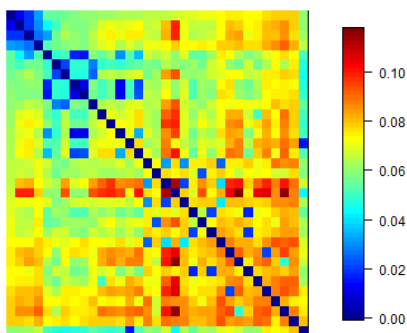
```
## g 0.003946461 0.003938269 0.003916016 0.003901434
## t 0.003952531 0.003975896 0.003939817 0.003894992
##
## $confidenceLevel
## [1] 0.95
##
## $lowerEndpointMatrix
##      a      c      g      t
## a 0.2464834 0.2466634 0.2427038 0.2385046
## c 0.2435866 0.2495944 0.2387691 0.2425160
## g 0.2461746 0.2451402 0.2423413 0.2405161
## t 0.2449741 0.2479176 0.2433798 0.2378003
##
## $upperEndpointMatrix
##      a      c      g      t
## a 0.2593811 0.2595657 0.2555034 0.2511945
## c 0.2563531 0.2625154 0.2514104 0.2552548
## g 0.2591573 0.2580959 0.2552239 0.2533507
## t 0.2579768 0.2609971 0.2563406 0.2506137
```

We fitted a first order markov model on all sequences. Our assumption in our simulated datasets is that in the sequence the occurrence of a nucleotide does not depend on the rest of the sequence. This violates the limited horizon: which is that the probability of being in a state at time t depends only on the state at time t minus 1. We used sample `{base}` function, which obviously samples without taking into account past states.

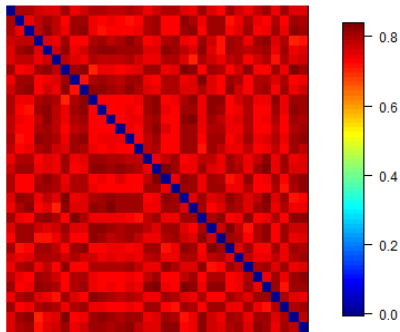
Question 2.3

To align the sequences for each dataset (the original dataset *lizards_sequences*, the first artificial dataset *artificial_dataset_1_1* and the second artificial dataset *artificial_dataset_1_2*), the *plsgenomics* package was used. The *.fasta*-files for the datasets were transformed to a *DNAStringSet* - class within R. The uncorrected distance matrices created represent the hamming distance between each of the sequences in each dataset. The results of these distance matrices are plotted as heatmaps (using *plsgenomics* package) :

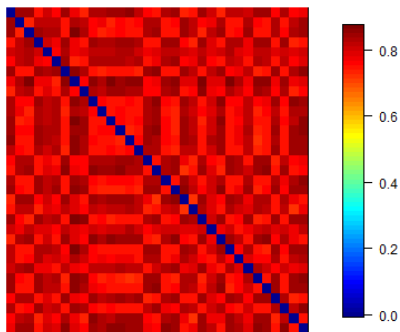
lizards_sequences



artificial_dataset_1_1



artificial_dataset_1_2



We see that for the original dataset, the alignment results are much better than for the artificial datasets. Based on the point that the artificial datasets were created by sampling randomly, the greater distances between the sequences compared to the distances within the original dataset make sense.

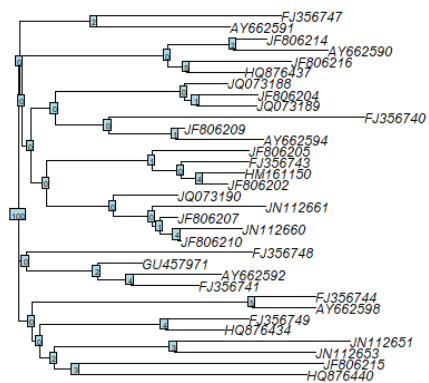
The R code for this Question 2.3 can be found in Appendix 2.3.

Question 3

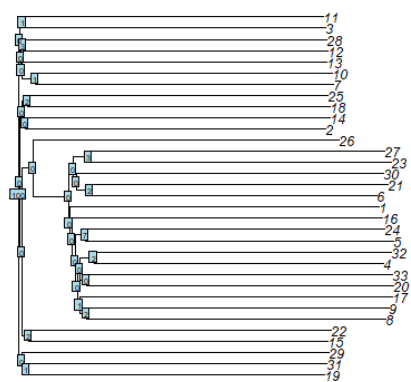
Question 3.1

Using the created distance matrix for each dataset (the original dataset *lizards_sequences*, the first artificial dataset *artificial_dataset_1_1* and the second artificial dataset *artificial_dataset_1_2*) with the aligned sequences, phylotrees were created. On top of that, a phylogenetic bootstrap analysis was performed. As a result, the bootstrap supports for the individual clades were integrated into the phylotrees.

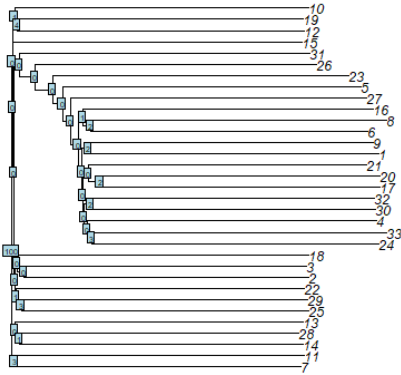
```
## detected function mkl_set_num_threads
lizards_sequences
```



artificial_dataset_1_1



artificial_dataset_1_2



The R code for the creation of the phylotrees and the bootstrap analysis can be found in Appendix 3.1.

Question 3.2

Different general characteristics can be compared between phylogenetic trees, e.g.:

- number of tips
- different tips
- number of nodes

On top of that, different quantitative distances can be calculated, e.g.:

- symmetric difference
- branch score

The distances can be only calculated if the tips are named equally. Since the artificial datasets (*artificial_dataset_1_1* and *artificial_dataset_1_2*) are not named as the original dataset (*lizard_sequences*), the distance measurements could be only processed for the comparison between the artificial datasets.

```
## => Comparing phylotree1 with phylotree2.
## Both trees have the same number of tips: 33.
## Tips in phylotree1 not in phylotree2 : JF806202, HM161150, FJ356743, JF806205, JQ073190, GU457971, F
## Tips in phylotree2 not in phylotree1 : 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18
## Both trees have the same number of nodes: 31.
## Both trees are unrooted.
## Both trees are not ultrametric.

## => Comparing phylotree1 with phylotree2.
## Both trees have the same number of tips: 33.
## Tips in phylotree1 not in phylotree2 : JF806202, HM161150, FJ356743, JF806205, JQ073190, GU457971, F
## Tips in phylotree2 not in phylotree1 : 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18
## Both trees have the same number of nodes: 31.
## Both trees are unrooted.
## Both trees are not ultrametric.

##      symmetric.difference    branch.score.difference
##              56.00000000              0.08985799
##      path.difference quadratic.path.difference
##              68.55654600              0.61762377
```

Appendix 1

Data Import of original dataset

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

Appendix 1.1

Reading and preparing original data

```
library(seqinr)
# reading original_dataset from fasta file
lizards_sequences = read.fasta("lizard_seqs.fasta")
```

Function code

```
library(seqinr)
get_artificial_sequence_dataset = function(original_dataset) {
  # creating empty variables which will be filled in following for-loop
  original_base_compositions = list()
  artificial_dataset = list()
  artificial_base_compositions = list()
  a_original = c(); c_original = c(); g_original = c(); t_original = c()
  a_artificial = c(); c_artificial = c(); g_artificial = c(); t_artificial = c()
  for (i in 1:length(original_dataset)) {
    # getting base compositions for each original sequence
    original_base_compositions[[i]] =
      seqinr::count(original_dataset[[i]],1)/length(original_dataset[[i]])
    # creating artificial sequences randomly drawn from the distribution
    # given by the base composition
    artificial_dataset[[as.character(i)]] = sample(x = c("a","c","g","t"),
                                                  size = length(original_dataset[[i]]),
                                                  rep = TRUE,
                                                  prob = original_base_compositions[[i]])

    # creating dataframe to compare base compositions
    # between original and artificial sequences
    artificial_base_compositions[[i]] =
      seqinr::count(artificial_dataset[[i]],1)/length(artificial_dataset[[i]])
  }
}
```



```

a_original = c(a_original, round(original_base_compositions[[i]][1],2))
a_artificial = c(a_artificial, round(artificial_base_compositions[[i]][1],2))
c_original = c(c_original, round(original_base_compositions[[i]][2],2))
c_artificial = c(c_artificial, round(artificial_base_compositions[[i]][2],2))
g_original = c(g_original, round(original_base_compositions[[i]][3],2))
g_artificial = c(g_artificial, round(artificial_base_compositions[[i]][3],2))
t_original = c(t_original, round(original_base_compositions[[i]][4],2))
t_artificial = c(t_artificial, round(artificial_base_compositions[[i]][4],2))
}
comparison_base_compositions = cbind(
  name_original = names(original_dataset), name_artificial = names(artificial_dataset),
  a_original, a_artificial, c_original, c_artificial,
  g_original, g_artificial, t_original, t_artificial
)
rownames(comparison_base_compositions) = 1:nrow(comparison_base_compositions)
print("comparison of base compositions
      between original and artificial datasets (values rounded): ")
print(comparison_base_compositions)
# saving fasta file
ape::write.dna(artificial_dataset, file = "artificial_dataset_1_1.fasta", format = "fasta",
               colsep = "")
}

```

Appendix 1.2

Replace the integers by letters

```

for (k in 1:33){
sequences_artificial[[k]][sequences_artificial[[k]] == 1] = "a"
sequences_artificial[[k]][sequences_artificial[[k]] == "2"] = "c"
sequences_artificial[[k]][sequences_artificial[[k]] == "3"] = "g"
sequences_artificial[[k]][sequences_artificial[[k]] == "4"] = "t"
}

```

Appendix 2

Appendix 2.3

```

library(seqinr)
library(DECIPHER)
library(plsgenomics)
library(ape)

# getting all datasets in DNASTringSet format

# original dataset
# readAAStringSet-function needs path of fasta file as input. The original
# dataset needs to be prepared and saved so that the fasta file does not
# include whitespaces anymore.
# reading original_dataset from fasta file
lizards_sequences = read.fasta("lizard_seqs.fasta")
# preparing data in fasta file (dna sequences include empty spaces which will be removed)
for (i in 1:length(lizards_sequences)) {

```

```

    lizards_sequences[[i]] = lizards_sequences[[i]][lizards_sequences[[i]] != " "]
  }
  # saving prepared fasta file
  ape::write.dna(lizards_sequences, file = "lizards_sequences_no_whitespaces.fasta",
                format = "fasta", colsep = "")
  # reading prepared fasta file as biostrings-object
  lizards_sequences = readDNAStringSet("lizards_sequences_no_whitespaces.fasta")

  # artificial_dataset_1_1
  artificial_dataset_1_1 = readDNAStringSet("artificial_dataset_1_1.fasta")

  # artificial_dataset_1_2
  artificial_dataset_1_2 = readDNAStringSet("artificial_dataset_1_2.fasta")

  # alligning sequences for each dataset
  sequence_alligning = function(dataset, name) {
    # alligning process
    sequences_aligned = AlignSeqs(dataset)
    # creating distance matrix
    dm_sequences_aligned = DistanceMatrix(sequences_aligned)
    # creating matrix heatmap
    heatmap_dm_sequences_aligned = matrix.heatmap(dm_sequences_aligned)
    dev.copy(png, paste("heatmap_", name, ".png", sep = ""))
    dev.off()
    return(sequences_aligned)
  }

  lizards_sequences_aligned = sequence_alligning(dataset = lizards_sequences,
                                                name = "lizards_sequences")
  artificial_dataset_1_1_aligned = sequence_alligning(artificial_dataset_1_1,
                                                    name = "artificial_dataset_1_1")
  artificial_dataset_1_2_aligned = sequence_alligning(artificial_dataset_1_2,
                                                    name = "artificial_dataset_1_2")

```

Appendix 3

Appendix 3.1

```

library(seqinr)
library(DECIPHER)
library(plsgenomics)
library(ape)

# creating phylotrees
create_phylotree = function(dataset_name) {
  distanceMatrix = readRDS(paste0("distanceMatrix_", dataset_name, ".RDS"))
  tree = nj(distanceMatrix)
  png(paste("phylotree_", dataset_name, ".png", sep = ""))
  plot(tree)
  dev.off()
  return(tree)
}

```

```

tree_lizards_sequences = create_phylotree("lizards_sequences")
tree_artificial_dataset_1_1 = create_phylotree("artificial_dataset_1_1")
tree_artificial_dataset_1_2 = create_phylotree("artificial_dataset_1_2")

# performing bootstrap analysis
bootstrap_analysis = function(dataset_name, tree_object) {
  distanceMatrix = readRDS(paste0("distanceMatrix_", dataset_name, ".RDS"))
  bootstrap_result = boot.phylo(phy = tree_object,
                                x = distanceMatrix,
                                FUN = function(x) {
                                  nj(x)
                                })
  png(paste("bootstrap_phylotree_", dataset_name, ".png", sep=""))
  plot(tree_object)
  nodelabels(bootstrap_result, cex=.6)
  dev.off()
}
bootstrap_analysis("lizards_sequences", tree_lizards_sequences)
bootstrap_analysis("artificial_dataset_1_1", tree_artificial_dataset_1_1)
bootstrap_analysis("artificial_dataset_1_2", tree_artificial_dataset_1_2)

```

Appendix 3.2

```

library(phangorn)
compare_phylotrees = function(phylogtree1, phylogtree2) {
  if(all(phylogtree1$tip.label == phylogtree2$tip.label)) {
    comparePhylo(phylogtree1, phylogtree2)
    treedist(phylogtree1, phylogtree2)
  } else {
    comparePhylo(phylogtree1, phylogtree2)
  }
}

# Comparing tree_lizards_sequences & tree_artificial_dataset_1_1
compare_phylotrees(tree_lizards_sequences, tree_artificial_dataset_1_1)
# Comparing tree_lizards_sequences & tree_artificial_dataset_1_2
compare_phylotrees(tree_lizards_sequences, tree_artificial_dataset_1_2)
# Comparing tree_artificial_dataset_1_1 & tree_artificial_dataset_1_2
compare_phylotrees(tree_artificial_dataset_1_1, tree_artificial_dataset_1_2)

```