

# Bioinformatics Lab2

*Milda Pocevičiute, Fanny Karelius, Rab Nawaz Jan Sher and Saman Zahid*

*24 November 2018*

## Question 1: DNA sequence acquisition and simulation

In this question an analysis of 3 nucleotide sequences is performed. One sequence is downloaded from GenBank database, and saved into fasta file. Here we load the data set:

```
lizards_sequences<-ape::read.FASTA("lizard_seqs.fasta")
```

### Q1.1

In this part we simulate DNA sequences with the same base composition and sequences length as the one from the GenBank:

```
len_list <- lengths(lizards_sequences)
new_seq_list <- list()
new_seq <- c()
for (i in 1:length(lizards_sequences)){
  # extract the base composition of each sequence
  freqx <- ape::base.freq(lizards_sequences[i])
  new_seq <- sample(c("a","c","g","t"),len_list[i],rep=TRUE,prob=freqx)
  new_seq_list[i] <- list(new_seq)
}

# convert result to the DNABin class
simulated_seq1 <- as.DNABin(new_seq_list)
names(simulated_seq1) <- names(lizards_sequences)
#ape::write.dna(simulated_seq1, file = "sim1_seq.fasta", format = "fasta", colsep = "")
```

## Base composition of the data from GenBank:

```
##           a           c           g           t
## 0.3121454 0.2052325 0.2307222 0.2518999
```

## Base composition of the simulated sequences:

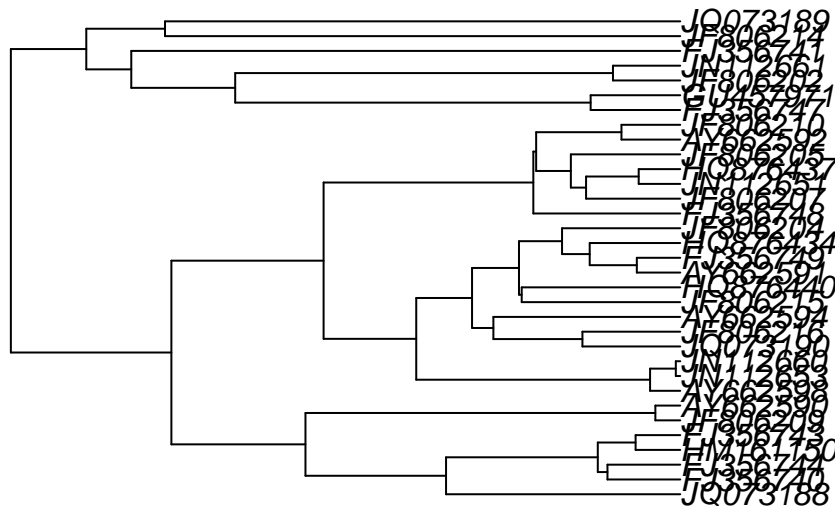
```
##           a           c           g           t
## 0.3100329 0.2065714 0.2311760 0.2522198
```

We can see that the original data, and the simulated data have quite similar base compositions.

### Q1.2

In this part we simulate another set of DNA data with similar characteristics as the original one. Now we use a simulated tree with 33 nodes from the original data and package phangorn to do this:

```
# Simulate a tree with 33 nodes
tree<- TreeSim::sim.bd.taxa(n=33, numbsim=1, lambda=1, mu=0)[[1]]
tree$tip.label <- names(lizards_sequences)
plot(tree, yaxt = "n")
```



We need to create a transition matrix, so we use the full dataset to compute that:

```
data <- c()
simulated_seq_tree <- list()

DNA_list <- unlist(as.character(lizards_sequences), use.names = FALSE)
DNA_unique <- c("a", "c", "g", "t")
matrix <- matrix(0, ncol = length(DNA_unique),
                 nrow = length(DNA_list))

for (i in 1:(length(DNA_list) - 1)) {
  index_of_i <- DNA_unique == DNA_list[i]
  index_of_i_plus_1 <- DNA_unique == DNA_list[i + 1]
  matrix[index_of_i, index_of_i_plus_1] = matrix[index_of_i, index_of_i_plus_1] + 1
}

# Transition matrix is finally computed
Q <- matrix / rowSums(matrix)
#Q_lower <- Q[lower.tri(Q)]
# Other parameters that are needed for phangorn package function
bf <- ape::base.freq(lizards_sequences)
size <- round(mean(len_list))

# Generating the data based on the tree

simulated_seq_tree <- phangorn::simSeq(
  tree,
```

```

    l = size,
    type = "DNA",
    bf = bf,
    Q = Q)

# change numbers into letters
simulated_seq_tree2 <- lapply(simulated_seq_tree,function(seq){
  #seq = a
  states = c("a","c","g","t")
  seq[which(seq==1)] <- states[1]
  seq[which(seq==2)] <- states[2]
  seq[which(seq==3)] <- states[3]
  seq[which(seq==4)] <- states[4]
  return(seq)
})

simulated_seq2 <- as.DNABin(simulated_seq_tree2)
# Write the simulated sequence to a fasta file
ape::write.dna(simulated_seq2, file = "sim2_seqs.fasta", format = "fasta", colsep = "")

## Base composition of the data from GenBank:
##          a          c          g          t
## 0.3121454 0.2052325 0.2307222 0.2518999

## Base composition of the simulated sequences Q1.1:
##          a          c          g          t
## 0.3100329 0.2065714 0.2311760 0.2522198

## Base composition of the simulated sequences Q1.2:
##          a          c          g          t
## 0.3143233 0.2019132 0.2302755 0.2534880

```

The base compositions of all simulated sequences are quite similar to the original one.

## Question 2: Sequence analysis

### Q2.1

```

## Individual base composition of sequence 1
## The data from the GenBank:
##          a          c          g          t
## 0.2898696 0.2026078 0.2437312 0.2637914
## The simulated data in Q1.1:
##          a          c          g          t
## 0.3266533 0.1673347 0.2454910 0.2605210
## The simulated data in Q1.2:
##          a          c          g          t
## 0.3333333 0.1981846 0.2249117 0.2435703
##
## Individual base composition of sequence 2
## The data from the GenBank:
##          a          c          g          t

```

```

## 0.3115541 0.2122554 0.2314507 0.2447398
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3196752 0.2122554 0.2270210 0.2410484
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3116490 0.1946546 0.2450832 0.2486132
##
## Individual base composition of sequence 3
## The data from the GenBank:
##      a      c      g      t
## 0.3130405 0.2099620 0.2348668 0.2421308
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3109651 0.2182636 0.2359045 0.2348668
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3161876 0.2007060 0.2344932 0.2486132
##
## Individual base composition of sequence 4
## The data from the GenBank:
##      a      c      g      t
## 0.2842942 0.2097416 0.2445328 0.2614314
## The simulated data in Q1.1:
##      a      c      g      t
## 0.2705649 0.2160555 0.2616452 0.2517344
## The simulated data in Q1.2:
##      a      c      g      t
## 0.2919818 0.2097832 0.2455875 0.2526475
##
## Individual base composition of sequence 5
## The data from the GenBank:
##      a      c      g      t
## 0.3059701 0.1987788 0.2360923 0.2591588
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3012212 0.1926730 0.2408412 0.2652646
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3040847 0.2017146 0.2365103 0.2576904
##
## Individual base composition of sequence 6
## The data from the GenBank:
##      a      c      g      t
## 0.2988084 0.2071494 0.2364803 0.2575619
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3107241 0.1989001 0.2263978 0.2639780
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3267776 0.1976803 0.2128089 0.2627332
##
## Individual base composition of sequence 7
## The data from the GenBank:

```

```

##          a          c          g          t
## 0.3137323 0.2068488 0.2338291 0.2455898
## The simulated data in Q1.1:
##          a          c          g          t
## 0.3099274 0.2096160 0.2262193 0.2542373
## The simulated data in Q1.2:
##          a          c          g          t
## 0.3126576 0.1896117 0.2269289 0.2708018
##
## Individual base composition of sequence 8
## The data from the GenBank:
##          a          c          g          t
## 0.2968127 0.2001992 0.2340637 0.2689243
## The simulated data in Q1.1:
##          a          c          g          t
## 0.2998008 0.2001992 0.2300797 0.2699203
## The simulated data in Q1.2:
##          a          c          g          t
## 0.3212305 0.2027231 0.2218860 0.2541604
##
## Individual base composition of sequence 9
## The data from the GenBank:
##          a          c          g          t
## 0.2998084 0.1992337 0.2346743 0.2662835
## The simulated data in Q1.1:
##          a          c          g          t
## 0.3026820 0.2030651 0.2327586 0.2614943
## The simulated data in Q1.2:
##          a          c          g          t
## 0.3323248 0.1921331 0.2244075 0.2511346
##
## Individual base composition of sequence 10
## The data from the GenBank:
##          a          c          g          t
## 0.3157534 0.2071918 0.2304795 0.2465753
## The simulated data in Q1.1:
##          a          c          g          t
## 0.3178082 0.2068493 0.2287671 0.2465753
## The simulated data in Q1.2:
##          a          c          g          t
## 0.3202219 0.1966717 0.2264246 0.2566818
##
## Individual base composition of sequence 11
## The data from the GenBank:
##          a          c          g          t
## 0.3121449 0.2041903 0.2325994 0.2510653
## The simulated data in Q1.1:
##          a          c          g          t
## 0.3093040 0.2080966 0.2407670 0.2418324
## The simulated data in Q1.2:
##          a          c          g          t
## 0.3101362 0.2032274 0.2355018 0.2511346
##
## Individual base composition of sequence 12

```

```

## The data from the GenBank:
##      a      c      g      t
## 0.3177408 0.2061677 0.2293832 0.2467082
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3088083 0.2082902 0.2424870 0.2404145
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3333333 0.2082703 0.2213817 0.2370146
##
## Individual base composition of sequence 13
## The data from the GenBank:
##      a      c      g      t
## 0.3187773 0.2034207 0.2241630 0.2536390
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3027656 0.2136099 0.2278020 0.2558224
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3212305 0.2002017 0.2264246 0.2521432
##
## Individual base composition of sequence 14
## The data from the GenBank:
##      a      c      g      t
## 0.3146384 0.2035273 0.2320988 0.2497354
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3130067 0.2076137 0.2266479 0.2527318
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3091276 0.2118003 0.2349975 0.2440746
##
## Individual base composition of sequence 15
## The data from the GenBank:
##      a      c      g      t
## 0.3196064 0.2033528 0.2255831 0.2514577
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3227802 0.1979622 0.2307132 0.2485444
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3131619 0.1966717 0.2370146 0.2531518
##
## Individual base composition of sequence 16
## The data from the GenBank:
##      a      c      g      t
## 0.2921236 0.2123629 0.2382851 0.2572283
## The simulated data in Q1.1:
##      a      c      g      t
## 0.2788845 0.2041833 0.2679283 0.2490040
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3126576 0.2097832 0.2309632 0.2465961
##

```

```

## Individual base composition of sequence 17
## The data from the GenBank:
##      a      c      g      t
## 0.2902903 0.2052052 0.2412412 0.2632633
## The simulated data in Q1.1:
##      a      c      g      t
## 0.301 0.206 0.213 0.280
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3111447 0.2037317 0.2249117 0.2602118
##
## Individual base composition of sequence 18
## The data from the GenBank:
##      a      c      g      t
## 0.3203540 0.2042478 0.2269027 0.2484956
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3118584 0.2053097 0.2251327 0.2576991
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3126576 0.2022189 0.2314675 0.2536561
##
## Individual base composition of sequence 19
## The data from the GenBank:
##      a      c      g      t
## 0.3170475 0.2003515 0.2274165 0.2551845
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3187083 0.2060372 0.2253422 0.2499122
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3217347 0.1936460 0.2344932 0.2501261
##
## Individual base composition of sequence 20
## The data from the GenBank:
##      a      c      g      t
## 0.2946429 0.2043651 0.2341270 0.2668651
## The simulated data in Q1.1:
##      a      c      g      t
## 0.2916667 0.2093254 0.2430556 0.2559524
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3227433 0.1996974 0.2269289 0.2506304
##
## Individual base composition of sequence 21
## The data from the GenBank:
##      a      c      g      t
## 0.3172554 0.1956522 0.2248641 0.2622283
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3125000 0.1915761 0.2289402 0.2669837
## The simulated data in Q1.2:
##      a      c      g      t
## 0.2980333 0.1956631 0.2486132 0.2576904

```

```

##
## Individual base composition of sequence 22
## The data from the GenBank:
##      a      c      g      t
## 0.3219225 0.2074689 0.2285615 0.2420470
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3301072 0.2132734 0.2336675 0.2229520
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3071104 0.2052446 0.2193646 0.2682804
##
## Individual base composition of sequence 23
## The data from the GenBank:
##      a      c      g      t
## 0.3144044 0.1966759 0.2250693 0.2638504
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3220456 0.2004147 0.1955771 0.2819627
## The simulated data in Q1.2:
##      a      c      g      t
## 0.2985376 0.2112960 0.2385275 0.2516389
##
## Individual base composition of sequence 24
## The data from the GenBank:
##      a      c      g      t
## 0.296 0.214 0.237 0.253
## The simulated data in Q1.1:
##      a      c      g      t
## 0.280 0.239 0.213 0.268
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3212305 0.2017146 0.2284418 0.2486132
##
## Individual base composition of sequence 25
## The data from the GenBank:
##      a      c      g      t
## 0.3175890 0.2054988 0.2284103 0.2485019
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3218188 0.2114910 0.2231230 0.2435671
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3267776 0.2052446 0.2218860 0.2460918
##
## Individual base composition of sequence 26
## The data from the GenBank:
##      a      c      g      t
## 0.3144453 0.2110934 0.2306464 0.2438148
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3220271 0.2138867 0.2282522 0.2358340
## The simulated data in Q1.2:
##      a      c      g      t

```



```

## 0.3197176 0.2022189 0.2284418 0.2496218
##
## Individual base composition of sequence 27
## The data from the GenBank:
##      a      c      g      t
## 0.3094527 0.2059701 0.2218905 0.2626866
## The simulated data in Q1.1:
##      a      c      g      t
## 0.2825871 0.1960199 0.2388060 0.2825871
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3116490 0.1951589 0.2365103 0.2566818
##
## Individual base composition of sequence 28
## The data from the GenBank:
##      a      c      g      t
## 0.3160083 0.2065142 0.2286902 0.2487872
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3092356 0.1947423 0.2317537 0.2642684
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3212305 0.1875946 0.2269289 0.2642461
##
## Individual base composition of sequence 29
## The data from the GenBank:
##      a      c      g      t
## 0.3207612 0.2020761 0.2280277 0.2491349
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3064684 0.2058111 0.2262193 0.2615012
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3091276 0.2012103 0.2309632 0.2586989
##
## Individual base composition of sequence 30
## The data from the GenBank:
##      a      c      g      t
## 0.2943396 0.2037736 0.2377358 0.2641509
## The simulated data in Q1.1:
##      a      c      g      t
## 0.2707547 0.1952830 0.2452830 0.2886792
## The simulated data in Q1.2:
##      a      c      g      t
## 0.3086233 0.2062532 0.2213817 0.2637418
##
## Individual base composition of sequence 31
## The data from the GenBank:
##      a      c      g      t
## 0.3200585 0.2071611 0.2232371 0.2495433
## The simulated data in Q1.1:
##      a      c      g      t
## 0.3104456 0.2209642 0.2264427 0.2421476
## The simulated data in Q1.2:

```

```
##           a           c           g           t
## 0.3040847 0.2178517 0.2329803 0.2450832
##
## Individual base composition of sequence 32
## The data from the GenBank:
##           a           c           g           t
## 0.2961117 0.2033898 0.2352941 0.2652044
## The simulated data in Q1.1:
##           a           c           g           t
## 0.2781655 0.2093719 0.2442672 0.2681954
## The simulated data in Q1.2:
##           a           c           g           t
## 0.3071104 0.2138174 0.2294503 0.2496218
##
## Individual base composition of sequence 33
## The data from the GenBank:
##           a           c           g           t
## 0.2932331 0.2030075 0.2406015 0.2631579
## The simulated data in Q1.1:
##           a           c           g           t
## 0.3243824 0.1750806 0.2491944 0.2513426
## The simulated data in Q1.2:
##           a           c           g           t
## 0.3010590 0.2067574 0.2324760 0.2597075
```

The CG contents are:

```
## In the data from the GenBank:
## [1] 0.4359547
## In the simulated data in Q1.1:
## [1] 0.4377474
## In the simulated data in Q1.2:
## [1] 0.4321888
```

A,C,G,T contents are:

```
# A,C,G,T content
cat("In the data from the GenBank:")

## In the data from the GenBank:
cat("\n")

print(ape::base.freq(lizards_sequences))

##           a           c           g           t
## 0.3121454 0.2052325 0.2307222 0.2518999
cat("In the simulated data in Q1.1:")

## In the simulated data in Q1.1:
cat("\n")

print(ape::base.freq(simulated_seq1))

##           a           c           g           t
```

```
## 0.3100329 0.2065714 0.2311760 0.2522198
```

```
cat("In the simulated data in Q1.2:")
```

```
## In the simulated data in Q1.2:
```

```
cat("\n")
```

```
print(ape::base.freq(simulated_seq2))
```

```
##           a           c           g           t  
## 0.3143233 0.2019132 0.2302755 0.2534880
```

The six versions of the nucleotide sequences of the lizard (all sequences starting from nucleotide 1, nucleotide 2 and nucleotide 3, as well as the reversed and complimented versions of them) were transformed into amino acid sequences using the EMBOSS Transeq function, and saved to fasta files:

```
# start from the second nucleotide  
#lizard <- as.character(lizards_sequences)  
#lizard_seq2 = lapply(lizard, function(a){return(as.character(a[-1]))})  
#lizard_seq3 = lapply(lizard, function(a){return(as.character(a[c(-1,-2)]))})  
#ape::write.dna(lizard_seq2, file = "lizard_seq2.fasta", format = "fasta", colsep = "")  
#ape::write.dna(lizard_seq3, file = "lizard_seq3.fasta", format = "fasta", colsep = "")
```

```
# Create 3 versions of lizard sequences
```

```
amino_lizard_seq <- read.fasta("lizard_amino_seqs.fasta")  
amino_lizard_seq2 <- read.fasta("amino_lizard2.fasta")  
amino_lizard_seq3 <- read.fasta("amino_lizard3.fasta")
```

```
amino_rev_com_lizard <- read.fasta("rev_com_lizard_amino_acid.fasta")  
amino_rev_com_lizard2 <- read.fasta("amino_reversed_lizard2.fasta")  
amino_rev_com_lizard3 <- read.fasta("amino_reversed_lizard3.fasta")
```

```
stop_code_count <- function(amino_acid_seq){  
  
  stop_count = c()  
  for (i in 1:length(amino_acid_seq)) {  
  
    sequence = amino_acid_seq[[i]]  
    stop <- sequence[which(sequence == "*")]  
    stop_count[i] = length(stop)  
  }  
  return(stop_count)  
}  
df = data.frame("lizard_seq" = sum(stop_code_count(amino_lizard_seq)),  
               "lizard_seq2" = sum(stop_code_count(amino_lizard_seq2)),  
               "lizard_seq3" = sum(stop_code_count(amino_lizard_seq3)),  
               "rev_comp_lizard_seq" = sum(stop_code_count(amino_rev_com_lizard)),  
               "rev_comp_lizard2" = sum(stop_code_count(amino_rev_com_lizard2)),  
               "rev_comp_lizard3" = sum(stop_code_count(amino_rev_com_lizard3))  
               )  
  
knitr::kable(df)
```

lizard_seq	lizard_seq2	lizard_seq3	rev_comp_lizard_seq	rev_comp_lizard2	rev_comp_lizard3
1269	645	1418	897	897	897

The above table shows the total number of stop codans in each of the fasta files of lizard sequences. We can see that the number of stop codans is very high in general (this is due to the fact that the sequences is yet not alligned, hence the same parts of the real RNA strand is repeated multiple times). However, when we started to translate into amino acids from the second nucleotide, we have considerably less stop codans than if we translated from the first or third nucleotide. The all reversed and complimented sequences have exactly same amount of stop codans. This is reasonable, as they all are exactly the same sequences apart from the last 1-3 nucleotides (the length of the sequences are different, and the first nucleotides became the last ones after this operation).

Therefore, it is more likely that the true sequence starts from nucleotide 2. However, to find out the true underlying sequence, we'd need to do much more work. We'd have to check all possible combinations of the different starts of all the sequences in the data set. Also probably we would need to query the databases and try to find if there is some information about the exact sequence.

## Q2.2

In this part we do markov chain analysis of the data sets. We combine all the sequences from each dataset to seperate lists, and use the *markovchainFit* function with bootstrapping to estimate the Markov models. This function fits the first order markov model to all three datasets.

```
lizard <- as.character(lizards_sequences)
sim1 <- as.character(simulated_seq1)
sim2 <- as.character(simulated_seq2)

flat_lizard <- unlist(lizard,use.names = FALSE)
remove_values = c("y","m","r","s")
flat_lizard <- flat_lizard[!flat_lizard %in% remove_values]
mc_lizard <- markovchainFit(flat_lizard,method = "bootstrap", nboot = 3,
                           name = "Bootstrap Mc",
                           parallel = TRUE)

# the original data contains some "weird nucleotides" - what to do with them?
which(flat_lizard == "y")

## integer(0)
which(flat_lizard == "m")

## integer(0)
which(flat_lizard == "r")

## integer(0)
which(flat_lizard == "s")

## integer(0)
flat_sim1 <- unlist(sim1,use.names = FALSE)
mc_sim1 <- markovchainFit(flat_sim1,method = "bootstrap", nboot = 3,
                          name = "Bootstrap Mc",
                          parallel = TRUE)
```

```

)

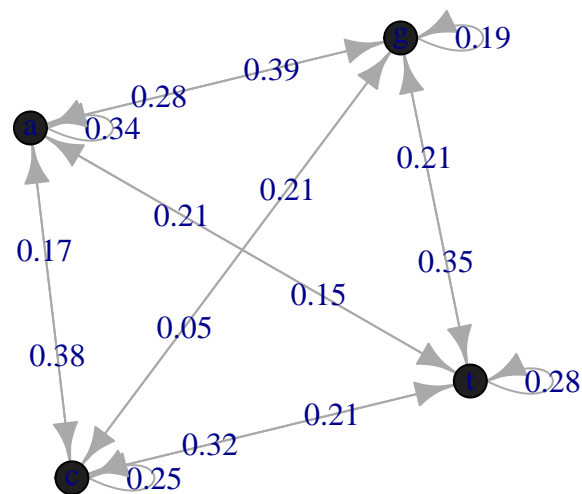
flat_sim2 <- unlist(sim2,use.names = FALSE)
mc_sim2 <- markovchainFit(flat_sim2,method = "bootstrap", nboot = 3,
                          name = "Bootstrap Mc",
                          parallel = TRUE
)

```

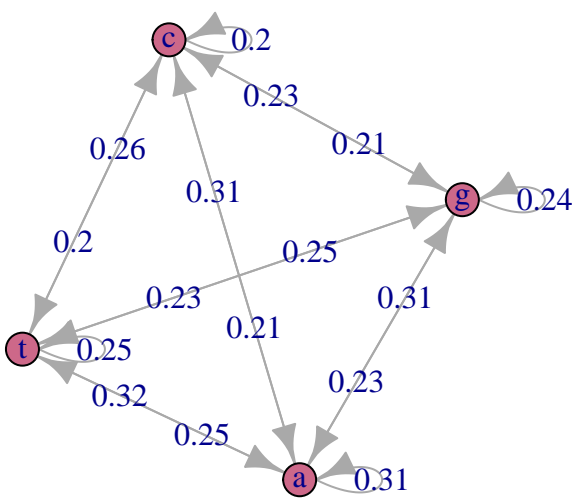
There are some additional letters in the lizard sequence. These are the nucleotides which we are not sure about where Y indicates unsurity in A or G while R shows ambiguity in C or T and M appears when there is no surity at all. We removed them before fitting the markov chains.

We see that the original data set from GenBank has a bit different structure than the simulated ones. Though the simulated also do not match in the structure. This is due to the fact, that the “true” underlying relations of the real lizard sequences have been not passed on to either of the simulations. Simulation 1 does not have any dependencies, while simulation 2 has some dependencies from the tree structure. However, these dependencies do not match the “real” dependencies in the real sequences.

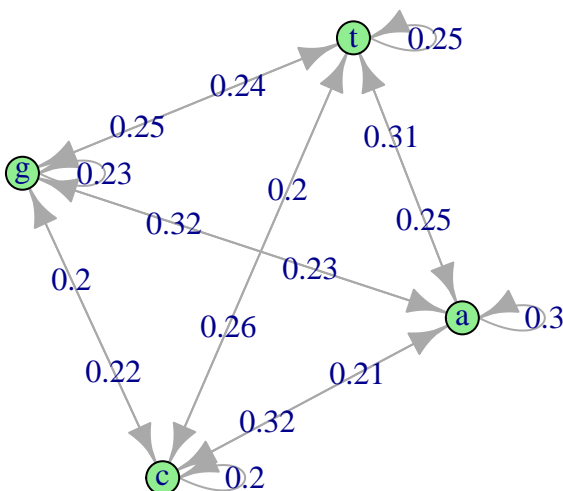
```
plot(mc_lizard$estimate)
```



```
plot(mc_sim1$estimate)
```



```
plot(mc_sim2$estimate)
```



In the first order Markov chains we assume that the nucleotide  $X_n$  only depends on a previous nucleotide, namely  $X_{n-1}$ . This is not a reasonable assumption, as DNA (or RNA) sequences contain the codons: the triplets of the nucleotides that determine what amino acids will be produced from it. Hence, it may be more reasonable that the most important dependences are between nucleotides in these triplets. In the lecture it was mentioned that there is some research supporting that a sixth order Markov Chain models perform better. That implies that the  $n$ th nucleotide depends on what were the previous 5 nucleotides.

### Q2.3

In this part the three DNA collections of the sequences are aligned using the *msa* package. Furthermore the distances between the alignments are calculated. They are used to produce the heatmaps

```

alignment_lizard <- msa("lizard_seqs.fasta", type="dna")

## use default substitution matrix
alignment_sim1 <- msa("sim1_seq.fasta", type="dna")

## use default substitution matrix
alignment_sim2 <- msa("sim2_seqs.fasta", type="dna")

## use default substitution matrix
alignment_lizard2 <- msaConvert(alignment_lizard, type="seqinr::alignment")
dist_lizard <- seqinr::dist.alignment(alignment_lizard2, matrix="identity")

alignment_sim1_2 <- msaConvert(alignment_sim1, type="seqinr::alignment")

```

```

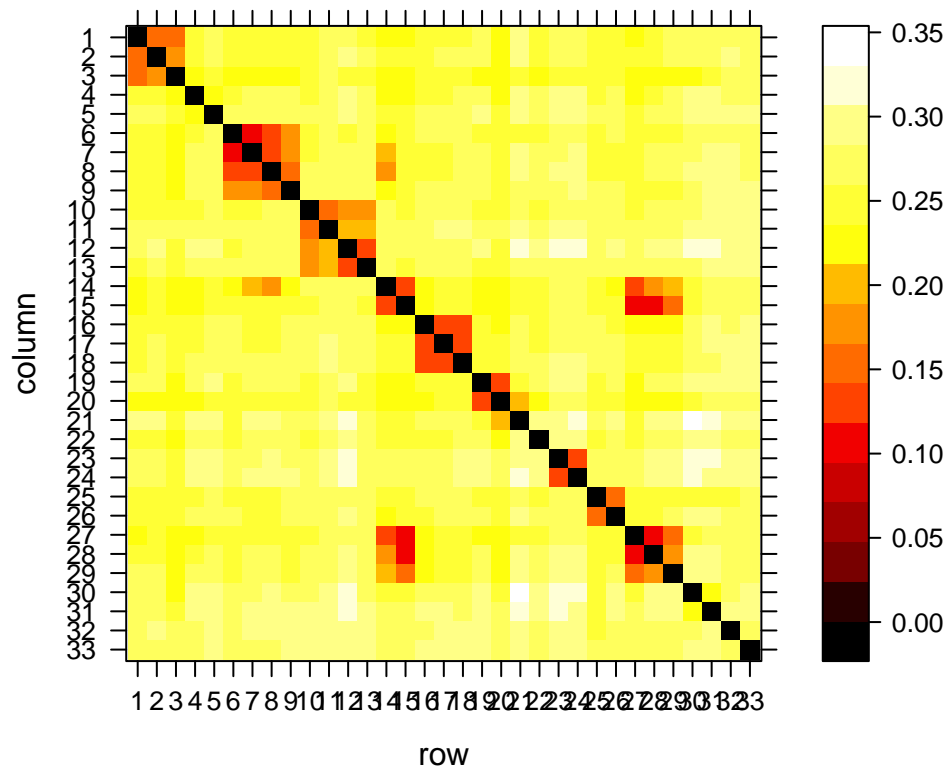
dist_sim1 <- seqinr::dist.alignment(alignment_sim1_2,matrix="identity")

alignment_sim2_2 <- msaConvert(alignment_sim2, type="seqinr::alignment")
dist_sim2 <- seqinr::dist.alignment(alignment_sim2_2,matrix="identity")

# the seqinr::dist.alignment computes the square root distance
dm_lizard <- as.matrix(dist_lizard)
dm_sim1 <- as.matrix(dist_sim1)
dm_sim2 <- as.matrix(dist_sim2)

# heatmaps
new.palette=colorRampPalette(c("black","red","yellow","white"),space="rgb")
levelplot(dm_lizard[1:ncol(dm_lizard),ncol(dm_lizard):1],col.regions=new.palette(20))

```

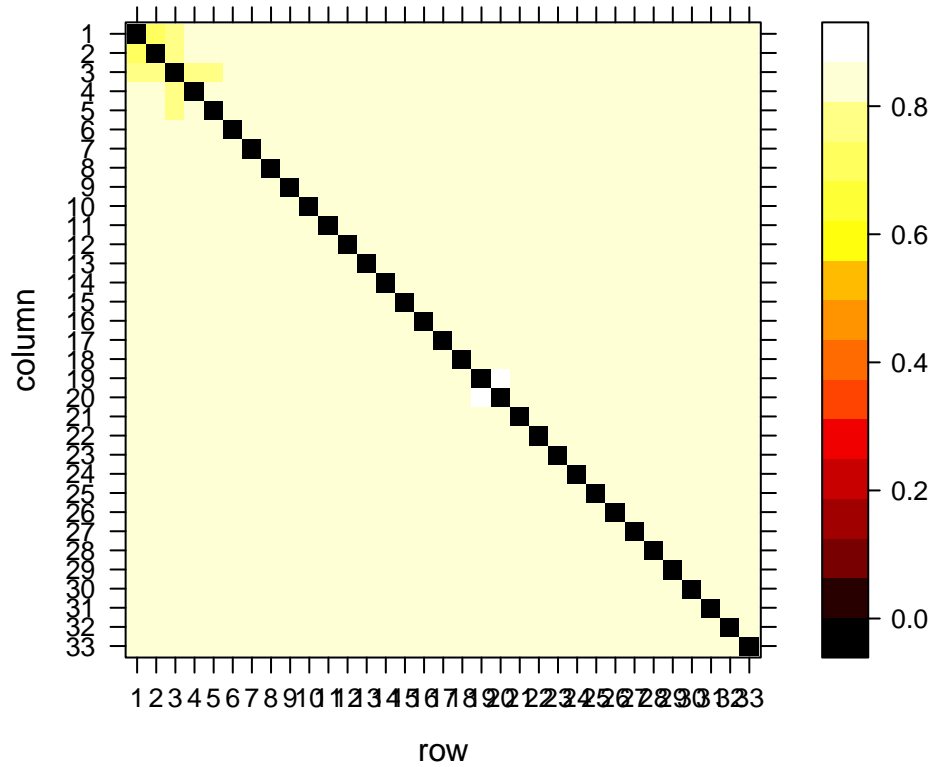


```

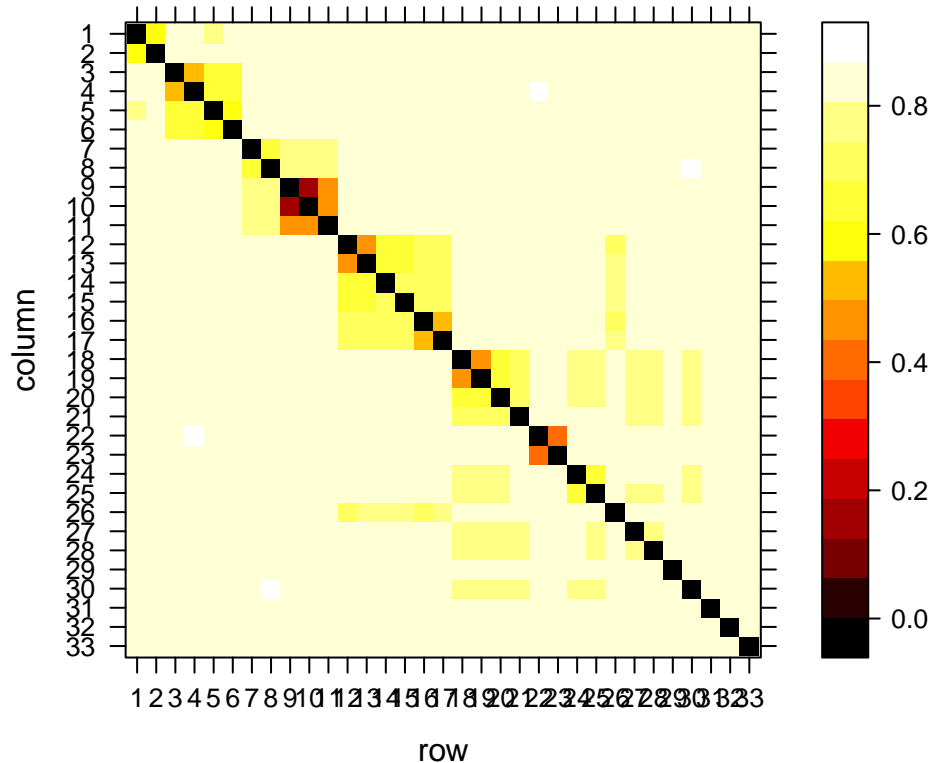
levelplot(dm_sim1[1:ncol(dm_sim1),ncol(dm_sim1):1],col.regions=new.palette(20))

```





```
levelplot(dm_sim2[1:ncol(dm_sim2),ncol(dm_sim2):1],col.regions=new.palette(20))
```



The plots illustrate how much overlaying the sequences have in each dataset. The smaller the value (the darker the colour), the more nucleotides are matching between the corresponding sequences of the DNA. Therefore, the sequences from the data from GenBank has much more overlaying parts then the two randomly generated sequences. This result is expected as sequencing of a DNA produces “fragments” of the real DNA sequence, and those “fragments” do contain the same bits of the original DNA. The artificially simulated DNA sequences from simulation 1 did not actually contain any structure (all nodes were assumed independent), but the simulated sequences from simulation 2 already contain some dependencies because of shared ancestors that were introduced by a tree structure (but this does not match the “true” dependencies from the original data). Hence they are way less likely to randomly contain long matching parts in their sequences.

### Question 3: Phylogeny reconstruction

#### Q3.1

```
library(phangorn)

upgma_tree <- function(x){
  tree <- upgma(x)
  tree$tip.label <- names(lizards_sequences)
  return(tree)
}
```

```

lizard_seq_tree <- upgma_tree(dm_lizard)
#plot(lizard_seq_tree, main="Phylogenetic Tree of Lizard DNA")

sim1_seq_tree <- upgma_tree(dm_sim1)
#plot(sim1_seq_tree, main="Phylogenetic Tree of Lizard DNA sim1")

sim2_seq_tree <- upgma_tree(dm_sim2)
#plot(sim2_seq_tree, main="Phylogenetic Tree of Lizard DNAs sim2")

boot_lizard <- boot.phylo(lizard_seq_tree, dm_lizard, FUN=upgma_tree, trees = TRUE, quiet = TRUE)
sim1_boot <- boot.phylo(sim1_seq_tree, dm_sim1, FUN=upgma_tree, trees = TRUE, quiet = TRUE)
sim2_boot <- boot.phylo(sim2_seq_tree, dm_sim2, FUN=upgma_tree, trees = TRUE, quiet = TRUE)

```

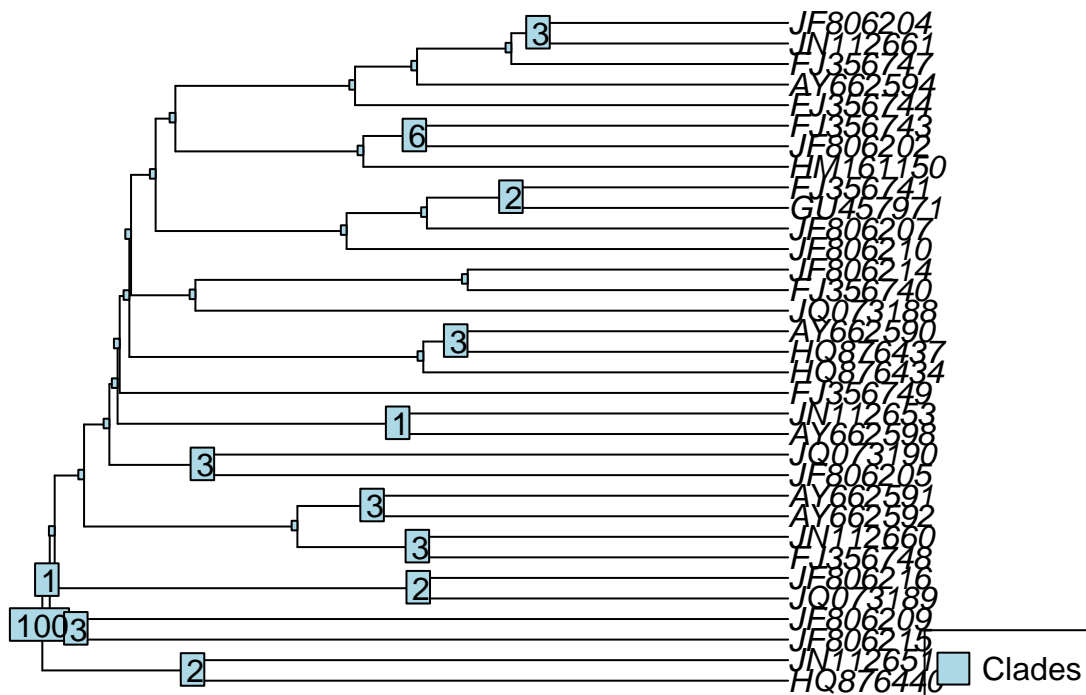
## Extracting Clades

```

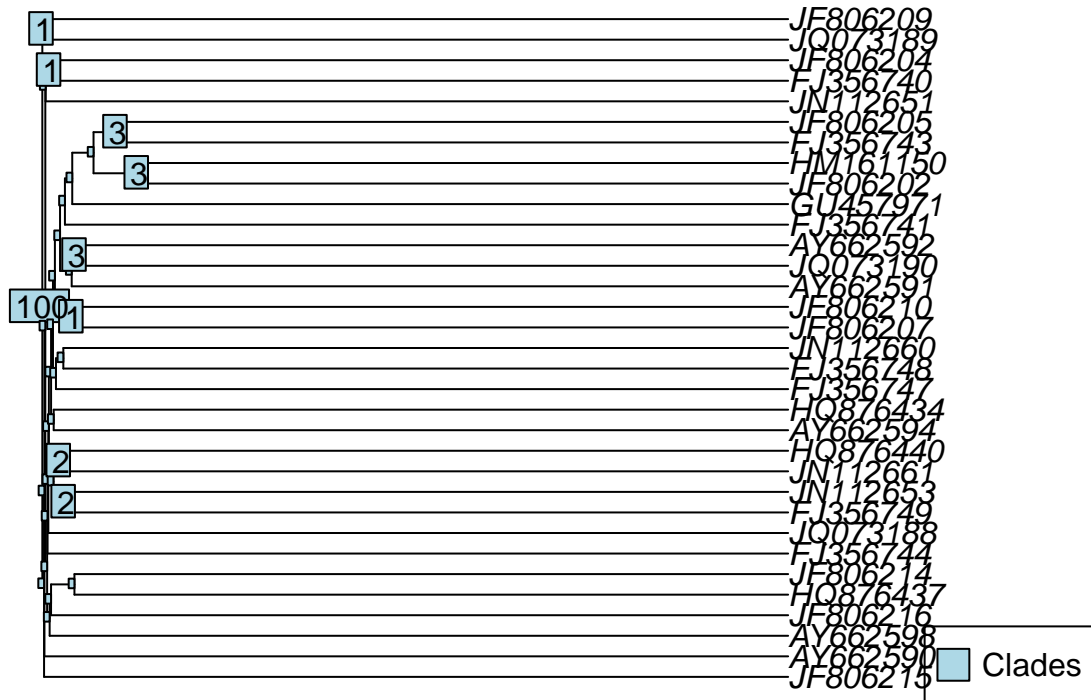
lizard_clade = prop.clades(lizard_seq_tree, boot_lizard$trees, rooted = TRUE)
sim1_clade = prop.clades(sim1_seq_tree, sim1_boot$trees, rooted = TRUE)
sim2_clade = prop.clades(sim2_seq_tree, sim2_boot$trees, rooted = TRUE)

```

## Phylogenetic Tree of Lizard DNA & Clades



# Sim 1 Phylogenetic Tree of Lizard DNA & Clades



A phylogenetic tree showing the relationships between 25 taxa. The taxa are listed on the right, grouped into three main clades: NY (top), JF (middle), and AY (bottom). Bootstrap values are indicated in blue boxes at the nodes. The root node has a bootstrap value of 100. The tree shows a clear hierarchical structure with several well-supported clades.

Clades

```
library(phangorn)

path_diff_sim1_lizard <- path.dist(sim1_seq_tree,lizard_seq_tree)
path_diff_sim2_lizard <- path.dist(sim2_seq_tree,lizard_seq_tree)
path_diff_sim1_sim2 <- path.dist(sim1_seq_tree,sim2_seq_tree)

felsenstein_diff_sim1_lizard <- KF.dist(sim1_seq_tree,lizard_seq_tree)
felsenstein_diff_sim2_lizard <- KF.dist(sim2_seq_tree,lizard_seq_tree)
felsenstein_diff_sim1_sim2 <- KF.dist(sim1_seq_tree,sim2_seq_tree)

result = data.frame(path_difference=c(path_diff_sim1_lizard,
                                     path_diff_sim2_lizard,
                                     path_diff_sim1_sim2),
                    felsenstein_distance=c(felsenstein_diff_sim1_lizard,
                                           felsenstein_diff_sim2_lizard,
                                           felsenstein_diff_sim1_sim2)
                    )

rownames(result) = c("sim1 vs lizard","sim2 vs lizard","sim1 vs sim2")
knitr::kable(result)
```

21

	path_difference	felsenstein_distance
sim2 vs lizard	105.46089	1.4103682
sim1 vs sim2	86.19745	0.8749121

## Geodesic Distance (Distance between edeges)

```
library(distory)
#distance between edeges using geodesic
geodesic_dist = dist.multiPhylo(list(lizard_seq_tree,sim1_seq_tree,sim2_seq_tree), method = "geodesic")
geodesic_dist = as.matrix(geodesic_dist)
geodesic_dist
```

```
##           1           2           3
## 1 0.000000 1.9198861 1.4501054
## 2 1.919886 0.0000000 0.8849977
## 3 1.450105 0.8849977 0.0000000
```

It can be seen from the path distance, Kuhner Felsensteins distance and edge distance that tree from simulation 1 differs more from expected lizard tree than simulation 2. It can also be observed that both simulations are closer to each other than is in terms of distance.

```
comp1 <- comparePhylo(lizard_seq_tree,sim1_seq_tree)
comp2 <- comparePhylo(lizard_seq_tree,sim2_seq_tree)
comp12 <- comparePhylo(sim1_seq_tree,sim2_seq_tree)
```

```
comp1$messages
```

```
## [1] "> Comparing lizard_seq_tree with sim1_seq_tree."
## [2] "Both trees have the same number of tips: 33."
## [3] "Both trees have the same tip labels."
## [4] "Both trees have the same number of nodes: 32."
## [5] "Both trees are rooted."
## [6] "Both trees are ultrametric."
## [7] "30 clades in lizard_seq_tree not in sim1_seq_tree."
## [8] "30 clades in sim1_seq_tree not in lizard_seq_tree."
## [9] "Branching times of clades in common between both trees: see ..$BT\n(node number in parentheses)"
```

```
comp2$messages
```

```
## [1] "> Comparing lizard_seq_tree with sim2_seq_tree."
## [2] "Both trees have the same number of tips: 33."
## [3] "Both trees have the same tip labels."
## [4] "Both trees have the same number of nodes: 32."
## [5] "Both trees are rooted."
## [6] "Both trees are ultrametric."
## [7] "30 clades in lizard_seq_tree not in sim2_seq_tree."
## [8] "30 clades in sim2_seq_tree not in lizard_seq_tree."
## [9] "Branching times of clades in common between both trees: see ..$BT\n(node number in parentheses)"
```

```
comp12$messages
```

```
## [1] "> Comparing sim1_seq_tree with sim2_seq_tree."
## [2] "Both trees have the same number of tips: 33."
## [3] "Both trees have the same tip labels."
## [4] "Both trees have the same number of nodes: 32."
```

```
## [5] "Both trees are rooted."
## [6] "Both trees are ultrametric."
## [7] "27 clades in sim1_seq_tree not in sim2_seq_tree."
## [8] "27 clades in sim2_seq_tree not in sim1_seq_tree."
## [9] "Branching times of clades in common between both trees: see ..$BT\n(node number in parentheses)"
```

The comparison in terms of tips, nodes, label and clades also implies that simulation 2 is more similar to lizard tree as it has 27 different clades from expected lizard tree while simulation 1 has 28 different clades

```
knitr::opts_chunk$set(echo = TRUE)

library(ape)
library(phangorn)
library(markovchain)
library(seqinr)
library(msa)
library(lattice)
lizards_sequences<-ape::read.FASTA("lizard_seqs.fasta")
len_list <- lengths(lizards_sequences)
new_seq_list <- list()
new_seq <- c()
for (i in 1:length(lizards_sequences)){
  # extract the base composition of each sequence
  freqx <- ape::base.freq(lizards_sequences[i])
  new_seq <- sample(c("a","c","g","t"),len_list[i],rep=TRUE,prob=freqx)
  new_seq_list[i] <- list(new_seq)
}

# convert result to the DNABin class
simulated_seq1 <- as.DNABin(new_seq_list)
names(simulated_seq1) <- names(lizards_sequences)
#ape::write.dna(simulated_seq1, file ="sim1_seq.fasta", format = "fasta", colsep = "")

cat("Base composition of the data from GenBank:")
cat("\n")
ape::base.freq(lizards_sequences)
cat("\n")
cat("Base composition of the simulated sequences:")
cat("\n")
ape::base.freq(simulated_seq1)

# Simulate a tree with 33 nodes
tree<- TreeSim::sim.bd.taxa(n=33, numbsim=1, lambda=1, mu=0)[[1]]
tree$tip.label <- names(lizards_sequences)
plot(tree, yaxt = "n")

data <- c()
simulated_seq_tree <- list()

DNA_list <- unlist(as.character(lizards_sequences),use.names = FALSE)
DNA_unique <- c("a","c","g","t")
matrix <- matrix(0, ncol = length(DNA_unique),
                 nrow = length(DNA_unique))

for (i in 1:(length(DNA_list) - 1)) {
```

```

    index_of_i <- DNA_unique == DNA_list[i]
    index_of_i_plus_1 <- DNA_unique == DNA_list[i + 1]
    matrix[index_of_i, index_of_i_plus_1] = matrix[index_of_i, index_of_i_plus_1] + 1
}

# Transition matrix is finally computed
Q <- matrix / rowSums(matrix)
#Q_lower <- Q[lower.tri(Q)]
# Other parameters that are needed for phangorn package function
bf <- ape::base.freq(lizards_sequences)
size <- round(mean(len_list))
# Generating the data based on the tree

simulated_seq_tree <- phangorn::simSeq(
  tree,
  l = size,
  type = "DNA",
  bf = bf,
  Q = Q)

# change numbers into letters
simulated_seq_tree2 <- lapply(simulated_seq_tree,function(seq){
  #seq = a
  states = c("a","c","g","t")
  seq[which(seq==1)] <- states[1]
  seq[which(seq==2)] <- states[2]
  seq[which(seq==3)] <- states[3]
  seq[which(seq==4)] <- states[4]
  return(seq)
})

simulated_seq2 <- as.DNABin(simulated_seq_tree2)
# Write the simulated sequence to a fasta file
ape::write.dna(simulated_seq2, file = "sim2_seqs.fasta", format = "fasta", colsep = "")

cat("Base composition of the data from GenBank:")
cat("\n")
ape::base.freq(lizards_sequences)
cat("\n")
cat("Base composition of the simulated sequences Q1.1:")
cat("\n")
ape::base.freq(simulated_seq1)
cat("\n")
cat("Base composition of the simulated sequences Q1.2:")
cat("\n")
ape::base.freq(simulated_seq2)

for (i in 1:length(lizards_sequences)){
  cat(paste0("Individual base composition of sequence ",i))
  cat("\n")
  cat("The data from the GenBank:")
  cat("\n")
  print(ape::base.freq(lizards_sequences[i]))
}

```



```

    cat("The simulated data in Q1.1:")
    cat("\n")
    print(ape::base.freq(simulated_seq1[i]))
    cat("The simulated data in Q1.2:")
    cat("\n")
    print(ape::base.freq(simulated_seq2[i]))
    cat("\n")
}

# CG content
comp_true <- ape::base.freq(lizards_sequences)
cg_true <- sum(comp_true[2:3])/sum(comp_true)
comp_s1 <- ape::base.freq(simulated_seq1)
cg_s1 <- sum(comp_s1[2:3])/sum(comp_s1)
comp_s2 <- ape::base.freq(simulated_seq2)
cg_s2 <- sum(comp_s2[2:3])/sum(comp_s2)

cat("In the data from the GenBank:")
cat("\n")
print(cg_true)
cat("In the simulated data in Q1.1:")
cat("\n")
print(cg_s1)
cat("In the simulated data in Q1.2:")
cat("\n")
print(cg_s2)
cat("\n")

# A,C,G,T content
cat("In the data from the GenBank:")
cat("\n")
print(ape::base.freq(lizards_sequences))
cat("In the simulated data in Q1.1:")
cat("\n")
print(ape::base.freq(simulated_seq1))
cat("In the simulated data in Q1.2:")
cat("\n")
print(ape::base.freq(simulated_seq2))

# start from the second nucleotide
#lizard <- as.character(lizards_sequences)
#lizard_seq2 = lapply(lizard, function(a){return(as.character(a[-1]))})
#lizard_seq3 = lapply(lizard, function(a){return(as.character(a[c(-1,-2)]))})
#ape::write.dna(lizard_seq2, file = "lizard_seq2.fasta", format = "fasta", colsep = "")
#ape::write.dna(lizard_seq3, file = "lizard_seq3.fasta", format = "fasta", colsep = "")

# Create 3 versions of lizard sequences

amino_lizard_seq <- read.fasta("lizard_amino_seqs.fasta")
amino_lizard_seq2 <- read.fasta("amino_lizard2.fasta")
amino_lizard_seq3 <- read.fasta("amino_lizard3.fasta")

amino_rev_com_lizard <- read.fasta("rev_com_lizard_amino_acid.fasta")
amino_rev_com_lizard2 <- read.fasta("amino_reversed_lizard2.fasta")
amino_rev_com_lizard3 <- read.fasta("amino_reversed_lizard3.fasta")

```

```

stop_code_count <- function(amino_acid_seq){

  stop_count = c()
  for (i in 1:length(amino_acid_seq)) {

    sequence = amino_acid_seq[[i]]
    stop <- sequence[which(sequence == "*")]
    stop_count[i] = length(stop)
  }
  return(stop_count)
}

df = data.frame("lizard_seq" = sum(stop_code_count(amino_lizard_seq)),
               "lizard_seq2" = sum(stop_code_count(amino_lizard_seq2)),
               "lizard_seq3" = sum(stop_code_count(amino_lizard_seq3)),
               "rev_comp_lizard_seq" = sum(stop_code_count(amino_rev_com_lizard)),
               "rev_comp_lizard2" = sum(stop_code_count(amino_rev_com_lizard2)),
               "rev_comp_lizard3" = sum(stop_code_count(amino_rev_com_lizard3))
               )

knitr::kable(df)

lizard <- as.character(lizards_sequences)
sim1 <- as.character(simulated_seq1)
sim2 <- as.character(simulated_seq2)

flat_lizard <- unlist(lizard,use.names = FALSE)
remove_values = c("y","m","r","s")
flat_lizard <- flat_lizard[!flat_lizard %in% remove_values]
mc_lizard <- markovchainFit(flat_lizard,method = "bootstrap", nboot = 3,
                           name = "Bootstrap Mc",
                           parallel = TRUE)

# the original data contains some "weird nucleotides" - what to do with them?
which(flat_lizard == "y")
which(flat_lizard == "m")
which(flat_lizard == "r")
which(flat_lizard == "s")

flat_sim1 <- unlist(sim1,use.names = FALSE)
mc_sim1 <- markovchainFit(flat_sim1,method = "bootstrap", nboot = 3,
                         name = "Bootstrap Mc",
                         parallel = TRUE
                         )

flat_sim2 <- unlist(sim2,use.names = FALSE)
mc_sim2 <- markovchainFit(flat_sim2,method = "bootstrap", nboot = 3,
                         name = "Bootstrap Mc",
                         parallel = TRUE
                         )

plot(mc_lizard$estimate)
plot(mc_sim1$estimate)

```

```

plot(mc_sim2$estimate)
alignment_lizard <- msa("lizard_seqs.fasta", type="dna")
alignment_sim1 <- msa("sim1_seq.fasta", type="dna")
alignment_sim2 <- msa("sim2_seqs.fasta", type="dna")

alignment_lizard2 <- msaConvert(alignment_lizard, type="seqinr::alignment")
dist_lizard <- seqinr::dist.alignment(alignment_lizard2,matrix="identity")

alignment_sim1_2 <- msaConvert(alignment_sim1, type="seqinr::alignment")
dist_sim1 <- seqinr::dist.alignment(alignment_sim1_2,matrix="identity")

alignment_sim2_2 <- msaConvert(alignment_sim2, type="seqinr::alignment")
dist_sim2 <- seqinr::dist.alignment(alignment_sim2_2,matrix="identity")

# the seqinr::dist.alignment computes the square root distance
dm_lizard <- as.matrix(dist_lizard)
dm_sim1 <- as.matrix(dist_sim1)
dm_sim2 <- as.matrix(dist_sim2)
# heatmaps
new.palette=colorRampPalette(c("black","red","yellow","white"),space="rgb")
levelplot(dm_lizard[1:ncol(dm_lizard),ncol(dm_lizard):1],col.regions=new.palette(20))
levelplot(dm_sim1[1:ncol(dm_sim1),ncol(dm_sim1):1],col.regions=new.palette(20))
levelplot(dm_sim2[1:ncol(dm_sim2),ncol(dm_sim2):1],col.regions=new.palette(20))

library(phangorn)

upgma_tree <- function(x){
  tree <- upgma(x)
  tree$tip.label <- names(lizards_sequences)
  return(tree)
}

lizard_seq_tree <- upgma_tree(dm_lizard)
#plot(lizard_seq_tree, main="Phylogenetic Tree of Lizard DNA")

sim1_seq_tree <- upgma_tree(dm_sim1)
#plot(sim1_seq_tree, main="Phylogenetic Tree of Lizard DNA sim1")

sim2_seq_tree <- upgma_tree(dm_sim2)
#plot(sim2_seq_tree, main="Phylogenetic Tree of Lizard DNAs sim2")

boot_lizard <- boot.phylo(lizard_seq_tree, dm_lizard, FUN=upgma_tree, trees = TRUE,quiet = TRUE)
sim1_boot <- boot.phylo(sim1_seq_tree, dm_sim1, FUN=upgma_tree, trees = TRUE,quiet = TRUE)
sim2_boot <- boot.phylo(sim2_seq_tree, dm_sim2, FUN=upgma_tree, trees = TRUE,quiet = TRUE)

lizard_clade = prop.clades(lizard_seq_tree,boot_lizard$trees,rooted = TRUE)

```

```

sim1_clade = prop.clades(sim1_seq_tree,sim1_boot$trees,rooted = TRUE)
sim2_clade = prop.clades(sim2_seq_tree,sim2_boot$trees,rooted = TRUE)
layout(1)
par(mar = rep(2, 4))
plot(lizard_seq_tree, main = "Phylogenetic Tree of Lizard DNA & Clades")
nodelabels(lizard_clade)
legend("bottomright", legend = c("Clades"), pch = 22,
      pt.bg = c("lightblue"), pt.cex = 2.5)

layout(1)
par(mar = rep(2, 4))
plot(sim1_seq_tree, main = "Sim 1 Phylogenetic Tree of Lizard DNA & Clades")
nodelabels(sim1_clade)
legend("bottomright", legend = c("Clades"), pch = 22,
      pt.bg = c("lightblue"), pt.cex = 2.5)

plot(sim2_seq_tree, main = "Sim 2 Phylogenetic Tree of Lizard DNA & Clades")
nodelabels(sim2_clade)
legend("bottomright", legend = c("Clades"), pch = 22,
      pt.bg = c("lightblue"), pt.cex = 2.5)

library(phangorn)
path_diff_sim1_lizard <- path.dist(sim1_seq_tree,lizard_seq_tree)
path_diff_sim2_lizard <- path.dist(sim2_seq_tree,lizard_seq_tree)
path_diff_sim1_sim2 <- path.dist(sim1_seq_tree,sim2_seq_tree)

felsenstein_diff_sim1_lizard <- KF.dist(sim1_seq_tree,lizard_seq_tree)
felsenstein_diff_sim2_lizard <- KF.dist(sim2_seq_tree,lizard_seq_tree)
felsenstein_diff_sim1_sim2 <- KF.dist(sim1_seq_tree,sim2_seq_tree)

result = data.frame(path_difference=c(path_diff_sim1_lizard,
                                     path_diff_sim2_lizard,
                                     path_diff_sim1_sim2),
                   felsenstein_distance=c(felsenstein_diff_sim1_lizard,
                                          felsenstein_diff_sim2_lizard,
                                          felsenstein_diff_sim1_sim2)
                   )
rownames(result) = c("sim1 vs lizard","sim2 vs lizard","sim1 vs sim2")
knitr::kable(result)
library(distory)
#distance between edeges using geodesic
geodesic_dist = dist.multiPhylo(list(lizard_seq_tree,sim1_seq_tree,sim2_seq_tree), method = "geodesic")
geodesic_dist = as.matrix(geodesic_dist)
geodesic_dist

comp1 <- comparePhylo(lizard_seq_tree,sim1_seq_tree)
comp2 <- comparePhylo(lizard_seq_tree,sim2_seq_tree)
comp12 <- comparePhylo(sim1_seq_tree,sim2_seq_tree)

```

```
comp1$messages  
comp2$messages  
comp12$messages
```