

Lab2 732A51 Bioinformatics Group 9

Duc Duong, Martin Smelik, Raymond Sseguya

2018 M11 27

Question 1: DNA sequence acquisition and simulation

1.1 Simulate an artificial DNA sequence dataset

Firstly, the code block below will get the lizards DNA from GenBank, and save it as a fasta file. (code provided by the teacher)

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

We create a function called `simulate_gene`, which take the lizards sequences as the input and simulate an AI gene base on the original sequences. When calling the function, AI gene is created automatically, and saved in a file called `AI_gene.fasta`. A message is returned to announce that the file is saved successfully.

```

#1.1
clean <- function(template_gene){
  nucleotide <- c("a", "t", "g", "c")
  for (i in 1:length(template_gene)) {
    #Remove the " " that created when reading a file
    template_gene[[i]] <- template_gene[[i]][template_gene[[i]]!= " "]
    #Remove the character that not nucleotide (eg: name of species...)
    template_gene[[i]] <- template_gene[[i]][match(template_gene[[i]], nucleotide)]
  }
  return(template_gene)
}

# Read and clean
lizards_sequences <- read.fasta("lizard_seqs.fasta")
lizards_sequences <- clean(lizards_sequences)

#Simulate AI gen function
simulate_gene <- function(template_gene)
{
  ai_gene <- list()
  gene_num <- length(template_gene)
  nucleotide <- c("a", "t", "g", "c")

  #Scan all gene and get some information
  for (i in 1:gene_num) {

    template_sequence <- template_gene[[i]]
    #get leng and base compotision of gene
    this_leng <- length(template_sequence)
    this_compotision = seqinr::count(template_sequence,1)/this_leng

    #generate a new sequence base on sample function
    this_sequence <- sample(nucleotide, size=this_leng ,prob = this_compotision, replace = TRUE)
    #print(this_sequence)

    #add to list
    ai_gene[i] <- list(this_sequence)
  }

  #write to a file
  ape::write.dna(ai_gene, file ="AI_gene.fasta", format = "fasta", colsep = " ")
  return("Created an AI gene and saved in file: AI_gene.fasta")
}

simulate_gene(lizards_sequences)

```

```
## [1] "Created an AI gene and saved in file: AI_gene.fasta"
```

Each sequence of the AI gene has the same base composition with the original gene. For example, here is the base composition of the first sequence:

```
ai_gene_1.1 <- read.fasta("AI_gene.fasta")
ai_gene_1.1 <- clean(ai_gene_1.1)

print("Base composition of the first sequence of the original gene: ")
```

```
## [1] "Base composition of the first sequence of the original gene: "
```

```
print(count(lizards_sequences[[1]],1)/length(lizards_sequences[[1]]))
```

```
##
##      a      c      g      t
## 0.2024048 0.5070140 0.0000000 0.2895792
```

```
print("Base composition of the first sequence of the AI gene: ")
```

```
## [1] "Base composition of the first sequence of the AI gene: "
```

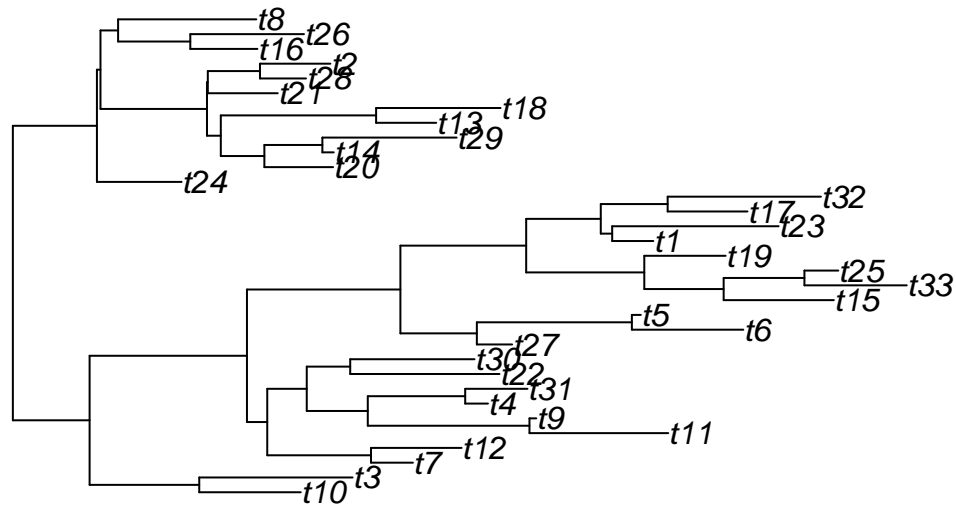
```
print(count(ai_gene_1.1[[1]],1)/length(ai_gene_1.1[[1]]))
```

```
##
##      a      c      g      t
## 0.256513 0.000000 0.000000 0.743487
```

1.2 Artificial DNA sequence dataset using phangorn::simSeq() function.

Here is the phylogenetic tree with 33 tips:

```
#1.2
tree <- rtree(length(lizards_sequences))
plot(tree)
```



Then, we create a transition Q matrix base of random number from 0.22 to 0.28. Which quite similar with the true number of the real data.

```
#the rates matrix
rates <- matrix(0, ncol = 4, nrow = 4)
rownames(rates) <- c("a", "c", "g", "t")
colnames(rates) <- c("a", "c", "g", "t")

#fill value to rates matrix
for (i in 1:4) {
  rate = runif(3, 0.22, 0.28)
  for (j in 1:4) {
    if (j==4)
    {
      rates[i,j]= 1- sum(rate)
    }else rates[i,j] = rate[j]
  }
}

#print the rates
rates
```

```
##           a           c           g           t
## a 0.2212496 0.2490569 0.2340297 0.2956638
## c 0.2246622 0.2726589 0.2660309 0.2366481
## g 0.2467697 0.2692747 0.2320388 0.2519168
```

```
## t 0.2673259 0.2796254 0.2758299 0.1772188
```

Finally, we create the second AI gene base on the `phangorn::simSeq()` function. We choose the length of all sequences equals 1000, Which more or less like the original sequence.

The second AI gene is saved as a fasta file with the name `AI_gene2.fasta`

```
#create the ai_gene 2
ai_gene_1.2 <- phangorn::simSeq(tree, l = 1000, Q=rates , type = "DNA")

#rename
for (i in 1:length(ai_gene_1.2)){
  ai_gene_1.2[[i]][ai_gene_1.2[[i]] == 1] = "a"
  ai_gene_1.2[[i]][ai_gene_1.2[[i]] == 2] = "c"
  ai_gene_1.2[[i]][ai_gene_1.2[[i]] == 3] = "g"
  ai_gene_1.2[[i]][ai_gene_1.2[[i]] == 4] = "t"
}

ape::write.dna(ai_gene_1.2, file ="AI_gene2.fasta", format = "fasta", colsep = "")
```

Question2: Sequence analysis

2.1: Report some basic statistics

Here is some basic statistics as the requirements:

```
#2.1
ai_gene_1.2 <- read.fasta("AI_gene2.fasta")
ai_gene_1.2 <- clean(ai_gene_1.2)

for (i in 1:length(lizards_sequences)) {
  cat(paste("For the sequences number: ", i , "\n"))
  print("The composition of lizards dataset:")
  print(round(count(lizards_sequences[[i]],2)/length(lizards_sequences[[i]]), 4))

  print("The composition of AI_gene1 dataset:")
  print(round(count(ai_gene_1.1[[i]],2)/length(ai_gene_1.1[[i]]), 4))

  print("The composition of Ai_gene2 dataset:")
  print(round(count(ai_gene_1.2[[i]],2)/length(ai_gene_1.2[[i]]), 4))

  cat("\n")
}
```

```
## For the sequences number: 1
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0501 0.0792 0.0000 0.0731 0.1032 0.2715 0.0000 0.1303 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0481 0.1553 0.0000 0.0862
## [1] "The composition of AI_gene1 dataset:"
##
```

```

##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0631 0.0000 0.0000 0.1934 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.1934 0.0000 0.0000 0.5491
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.999 0.000
##      ta      tc      tg      tt
## 0.000 0.000 0.000 0.000
##
## For the sequences number: 2
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.5814 0.0000 0.0000 0.1868 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.1868 0.0000 0.0000 0.0447
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.5925 0.1776 0.0000 0.0000 0.1776 0.0520 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.051 0.123 0.064 0.000 0.133 0.246 0.125
##      ta      tc      tg      tt
## 0.000 0.055 0.134 0.068
##
## For the sequences number: 3
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.4825 0.2041 0.0000 0.0000 0.2041 0.1090 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.9997 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.250 0.117 0.000 0.130 0.122 0.057 0.000 0.061 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.125 0.066 0.000 0.071
##
## For the sequences number: 4
## [1] "The composition of lizards dataset:"
##

```

```

##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0525 0.0833 0.0714 0.0000 0.1021 0.2706 0.1298 0.0000 0.0525 0.1487
##      gg      gt      ta      tc      tg      tt
## 0.0823 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.2329 0.0000 0.0000 0.2319 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.2309 0.0000 0.0000 0.3033
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.250 0.128 0.110 0.000 0.134 0.086 0.062
##      ta      tc      tg      tt
## 0.000 0.104 0.068 0.057
##
## For the sequences number: 5
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0482 0.1879 0.0000 0.0000 0.1879 0.5753 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.9993
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.068 0.073 0.124 0.000 0.061 0.067 0.126 0.000 0.136 0.114 0.230 0.000
##      ta      tc      tg      tt
## 0.000 0.000 0.000 0.000
##
## For the sequences number: 6
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0550 0.0000 0.0678 0.0843 0.0000 0.0000 0.0000 0.0000 0.0532 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0678 0.1366 0.0990 0.0000 0.1210 0.3144
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0522 0.1247 0.0467 0.0000 0.1100 0.2915 0.1292 0.0000 0.0614 0.1146
##      gg      gt      ta      tc      tg      tt
## 0.0687 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.062 0.115 0.058 0.000 0.119 0.266 0.128

```

```

##      ta      tc      tg      tt
## 0.000 0.055 0.131 0.065
##
## For the sequences number: 7
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.4815 0.2044 0.0000 0.0000 0.2044 0.1093 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.9997 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.059 0.062 0.000 0.130 0.061 0.049 0.000 0.131 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.130 0.131 0.000 0.246
##
## For the sequences number: 8
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.2679 0.2351 0.0000 0.0000 0.2341
##      gg      gt      ta      tc      tg      tt
## 0.2620 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.2420 0.2629 0.0000 0.0000 0.2620 0.2321
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.238 0.000 0.120 0.136 0.000 0.000 0.000 0.000 0.129 0.000 0.055 0.062
##      ta      tc      tg      tt
## 0.126 0.000 0.071 0.062
##
## For the sequences number: 9
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.4895 0.2107 0.0000 0.0000 0.2098 0.0891
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.999 0.000 0.000 0.000 0.000 0.000 0.000

```



```

##      ta      tc      tg      tt
## 0.000 0.000 0.000 0.000
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.238 0.000 0.247 0.000 0.000 0.000 0.000 0.000 0.246 0.000 0.268 0.000
##      ta      tc      tg      tt
## 0.000 0.000 0.000 0.000
##
## For the sequences number: 10
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.4774 0.2065 0.0000 0.0000 0.2065 0.1092 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.4438 0.0000 0.0000 0.2267 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.2264 0.0000 0.0000 0.1027
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.067 0.000 0.171 0.000 0.000 0.000 0.000 0.000 0.170 0.000 0.591 0.000
##      ta      tc      tg      tt
## 0.000 0.000 0.000 0.000
##
## For the sequences number: 11
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0508 0.0653 0.0771 0.0110 0.0533 0.0717 0.0384 0.0877 0.0540 0.0632
##      gg      gt      ta      tc      tg      tt
## 0.1083 0.0863 0.0462 0.0508 0.0884 0.0472
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0504 0.0611 0.1268 0.0000 0.0611
##      gg      gt      ta      tc      tg      tt
## 0.0618 0.1286 0.0000 0.1268 0.1282 0.2550
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.056 0.049 0.124 0.000 0.055 0.063 0.129 0.000 0.118 0.136 0.269 0.000
##      ta      tc      tg      tt
## 0.000 0.000 0.000 0.000
##
## For the sequences number: 12
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.4722 0.2059 0.0000 0.0000 0.2062 0.1092 0.0000 0.0000 0.0000 0.0000

```

```

##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.1067 0.0000 0.0000 0.2169 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.2169 0.0000 0.0000 0.4591
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.267 0.000 0.253 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.000 0.253 0.000 0.226
##
## For the sequences number: 13
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.6426 0.0000 0.1536 0.0000 0.0000 0.0000 0.0000 0.0000 0.1536 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0499 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.9996 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.080 0.132 0.059 0.000 0.137 0.242 0.117
##      ta      tc      tg      tt
## 0.000 0.055 0.121 0.056
##
## For the sequences number: 14
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.4780 0.2062 0.0000 0.0000 0.2062 0.1079 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.9996 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.052 0.115 0.070 0.000 0.118 0.236 0.134 0.000 0.067 0.137 0.070 0.000
##      ta      tc      tg      tt
## 0.000 0.000 0.000 0.000

```

```

##
## For the sequences number: 15
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.6419 0.0000 0.1521 0.0000 0.0000 0.0000 0.0000 0.0000 0.1517 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0509 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.9996 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.217 0.130 0.000 0.120 0.131 0.077 0.000 0.071 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.118 0.072 0.000 0.063
##
## For the sequences number: 16
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.2709 0.0956 0.1365 0.0000 0.1454
##      gg      gt      ta      tc      tg      tt
## 0.0418 0.0498 0.0000 0.0857 0.1006 0.0707
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.2201 0.2510 0.0000 0.0000 0.2510 0.2769
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.265 0.234 0.000 0.000 0.234 0.266 0.000
##      ta      tc      tg      tt
## 0.000 0.000 0.000 0.000
##
## For the sequences number: 17
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.053 0.081 0.000 0.071 0.095 0.272 0.000 0.135 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.056 0.150 0.000 0.084
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.091 0.000 0.217 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.000 0.217 0.000 0.474

```

```

## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.069 0.063 0.127 0.000 0.065 0.060 0.117
##      ta      tc      tg      tt
## 0.000 0.124 0.119 0.255
##
## For the sequences number: 18
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.2411 0.1179 0.0935 0.0000 0.1165 0.1112 0.0927 0.0000 0.0949 0.0913
##      gg      gt      ta      tc      tg      tt
## 0.0407 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0552 0.0000 0.1023 0.0704 0.0000 0.0000 0.0000 0.0000 0.0973 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.1979 0.1487 0.0750 0.0000 0.1437 0.1090
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.048 0.000 0.000 0.183 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.182 0.000 0.000 0.586
##
## For the sequences number: 19
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.4770 0.2043 0.0000 0.0000 0.2036 0.1120 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.1042 0.0000 0.0000 0.2211 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.2208 0.0000 0.0000 0.4535
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.068 0.000 0.129 0.076 0.000 0.000 0.000 0.000 0.143 0.000 0.222 0.115
##      ta      tc      tg      tt
## 0.061 0.000 0.130 0.055
##
## For the sequences number: 20
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0536 0.0784 0.0724 0.0000 0.0982 0.2679 0.1349 0.0000 0.0526 0.1538
##      gg      gt      ta      tc      tg      tt
## 0.0873 0.0000 0.0000 0.0000 0.0000 0.0000

```

```

## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0982 0.2093 0.0000 0.0000 0.2103 0.4812
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.066 0.061 0.069 0.057 0.066 0.054 0.070 0.072 0.063 0.059 0.055 0.061
##      ta      tc      tg      tt
## 0.058 0.088 0.044 0.056
##
## For the sequences number: 21
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.9993 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.9993
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.552 0.185
##      ta      tc      tg      tt
## 0.000 0.000 0.186 0.076
##
## For the sequences number: 22
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.4691 0.2084 0.0000 0.0000 0.2081 0.1134 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.9997 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.057 0.181 0.000 0.000 0.180 0.581 0.000 0.000 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.000 0.000 0.000 0.000
##
## For the sequences number: 23

```

```

## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0435 0.1811 0.0000 0.0000 0.1811 0.5895 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.9993
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.082 0.073 0.058 0.061 0.070 0.069 0.054 0.062 0.059 0.052 0.059 0.055
##      ta      tc      tg      tt
## 0.063 0.061 0.054 0.067
##
## For the sequences number: 24
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.053 0.082 0.000 0.079 0.107 0.253 0.000 0.129 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.054 0.154 0.000 0.088
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.241 0.000 0.000 0.259 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.258 0.000 0.000 0.241
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.082 0.074 0.060 0.062 0.071 0.056 0.068 0.063 0.060 0.063 0.052 0.052
##      ta      tc      tg      tt
## 0.065 0.066 0.047 0.058
##
## For the sequences number: 25
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.4748 0.2073 0.0000 0.0000 0.2073 0.1103 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.9996
## [1] "The composition of Ai_gene2 dataset:"
##

```

```

##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.081 0.000 0.118 0.065 0.000 0.000 0.000 0.000 0.119 0.000 0.263 0.114
##      ta      tc      tg      tt
## 0.064 0.000 0.115 0.060
##
## For the sequences number: 26
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.6301 0.0000 0.1584 0.0000 0.0000 0.0000 0.0000 0.0000 0.1584 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0527 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.9996 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.251 0.000 0.259 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.000 0.259 0.000 0.230
##
## For the sequences number: 27
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.5453 0.0000 0.1910 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.1910 0.0000 0.0716
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.000 0.000 0.000 0.999
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.252 0.121 0.131 0.000 0.127 0.070 0.049
##      ta      tc      tg      tt
## 0.000 0.126 0.054 0.069
##
## For the sequences number: 28
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.4767 0.2044 0.0000 0.0000 0.2048 0.1103 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##

```

```

##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.4476 0.0000 0.0000 0.2217 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.2217 0.0000 0.0000 0.1086
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.249 0.113 0.140 0.000 0.109 0.051 0.059
##      ta      tc      tg      tt
## 0.000 0.144 0.055 0.079
##
## For the sequences number: 29
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.4725 0.2058 0.0000 0.0000 0.2058 0.1148 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.9997 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.073 0.060 0.000 0.129 0.065 0.063 0.000 0.123 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.124 0.128 0.000 0.234
##
## For the sequences number: 30
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.2406 0.2274 0.0000 0.0000 0.2264 0.3047
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.9991 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.056 0.128 0.000 0.060 0.119 0.245 0.000 0.136 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.069 0.127 0.000 0.059
##
## For the sequences number: 31
## [1] "The composition of lizards dataset:"
##

```



```

##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0698 0.1278 0.0000 0.0511 0.1129 0.3265 0.0000 0.1037 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0668 0.0888 0.0000 0.0515
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.9996
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.060 0.000 0.000 0.192 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.192 0.000 0.000 0.555
##
## For the sequences number: 32
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0489 0.0818 0.0000 0.0728 0.0977 0.2702 0.0000 0.1326 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0568 0.1476 0.0000 0.0907
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0867 0.1944 0.0000 0.0000 0.1934 0.5244 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.077 0.130 0.000 0.060 0.118 0.227 0.000 0.126 0.000 0.000 0.000 0.000
##      ta      tc      tg      tt
## 0.072 0.114 0.000 0.075
##
## For the sequences number: 33
## [1] "The composition of lizards dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0537 0.0730 0.0763 0.0000 0.0945 0.2760 0.1332 0.0000 0.0548 0.1536
##      gg      gt      ta      tc      tg      tt
## 0.0838 0.0000 0.0000 0.0000 0.0000 0.0000
## [1] "The composition of AI_gene1 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc
## 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
##      gg      gt      ta      tc      tg      tt
## 0.0354 0.1515 0.0000 0.0000 0.1525 0.6595
## [1] "The composition of Ai_gene2 dataset:"
##
##      aa      ac      ag      at      ca      cc      cg      ct      ga      gc      gg      gt
## 0.000 0.000 0.000 0.000 0.000 0.066 0.119 0.066 0.000 0.126 0.246 0.129

```

```
##      ta      tc      tg      tt
## 0.000 0.059 0.136 0.052
```

2.2: Markov chain

We decided to use the `markovchain` library to fit the markovchain model. Here is the result:

```
library(markovchain)
markovchainFit(lizards_sequences)

## $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 0.6671729 0.2394909 0.06098894 0.03234722
## c 0.3756621 0.5070469 0.05169315 0.06559781
## g 0.2987863 0.2020121 0.35499840 0.14420313
## t 0.2284000 0.2542000 0.18820000 0.32920000
##
##
## $standardError
##      a      c      g      t
## a 0.004501558 0.002697045 0.001361035 0.0009912021
## c 0.004215073 0.004897004 0.001563590 0.0017613720
## g 0.006907547 0.005679789 0.007529336 0.0047987797
## t 0.006758698 0.007130217 0.006135145 0.0081141851
##
## $confidenceLevel
## [1] 0.95
##
## $lowerEndpointMatrix
##      a      c      g      t
## a 0.6597685 0.2350547 0.05875024 0.03071684
## c 0.3687289 0.4989921 0.04912127 0.06270061
## g 0.2874244 0.1926697 0.34261375 0.13630984
## t 0.2172829 0.2424718 0.17810859 0.31585335
##
## $upperEndpointMatrix
##      a      c      g      t
## a 0.6745773 0.2439272 0.06322765 0.03397761
## c 0.3825953 0.5151018 0.05426503 0.06849500
## g 0.3101482 0.2113546 0.36738306 0.15209642
## t 0.2395171 0.2659282 0.19829141 0.34254665

markovchainFit(ai_gene_1.1)

## $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 0.84615899 0.02716809 0.01137580 0.1152971
## c 0.15893643 0.66373451 0.06257497 0.1147541
## g 0.05310559 0.04611801 0.64456522 0.2562112
## t 0.14352775 0.02384513 0.06796278 0.7646643
##
##
## $standardError
##           a           c           g           t
## a 0.005320808 0.0009534132 0.0006169395 0.001964087
## c 0.005636894 0.0115192805 0.0035369464 0.004789744
## g 0.002871621 0.0026760385 0.0100043874 0.006307483
## t 0.002441810 0.0009952765 0.0016802716 0.005636108
##
## $confidenceLevel
## [1] 0.95
##
## $lowerEndpointMatrix
##           a           c           g           t
## a 0.83740704 0.02559987 0.01036103 0.1120665
## c 0.14966456 0.64478698 0.05675721 0.1068757
## g 0.04838219 0.04171632 0.62810946 0.2458363
## t 0.13951133 0.02220805 0.06519898 0.7553938
##
## $upperEndpointMatrix
##           a           c           g           t
## a 0.85491094 0.02873632 0.01239058 0.1185277
## c 0.16820829 0.68268204 0.06839273 0.1226325
## g 0.05782899 0.05051970 0.66102097 0.2665861
## t 0.14754417 0.02548222 0.07072658 0.7739349
```

```
markovchainFit(ai_gene_1.2)
```

```
## $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 0.33255554 0.2088681 0.2150358 0.2435406
## c 0.1408242 0.4219035 0.1968352 0.2404372
## g 0.1350098 0.1798413 0.5377718 0.1473771
## t 0.1708525 0.2460795 0.1727391 0.4103290
##
##
## $standardError
##           a           c           g           t
## a 0.007445478 0.005900607 0.005987093 0.006371567
## c 0.004003986 0.006930433 0.004733748 0.005231842
## g 0.003730179 0.004305184 0.007444680 0.003897283
## t 0.004488354 0.005386593 0.004513067 0.006955727
##
## $confidenceLevel
```

```
## [1] 0.95
##
## $lowerEndpointMatrix
##      a      c      g      t
## a 0.3203087 0.1991625 0.2051879 0.2330603
## c 0.1342383 0.4105039 0.1890488 0.2318315
## g 0.1288742 0.1727599 0.5255264 0.1409666
## t 0.1634698 0.2372193 0.1653157 0.3988878
##
## $upperEndpointMatrix
##      a      c      g      t
## a 0.3448021 0.2185738 0.2248837 0.2540209
## c 0.1474102 0.4333030 0.2046215 0.2490428
## g 0.1411454 0.1869227 0.5500172 0.1537876
## t 0.1782352 0.2549396 0.1801624 0.4217701
```

Markov chain order: I'll do it later. There is an error in here

```
#2.2
#The fitHigherOrder function work only with a list, not list of list
#So, my idea is take some random sequences, see the order and then make the conclusion
fitHigherOrder(lizards_sequences[[1]])
```

2.3: Align the sequences

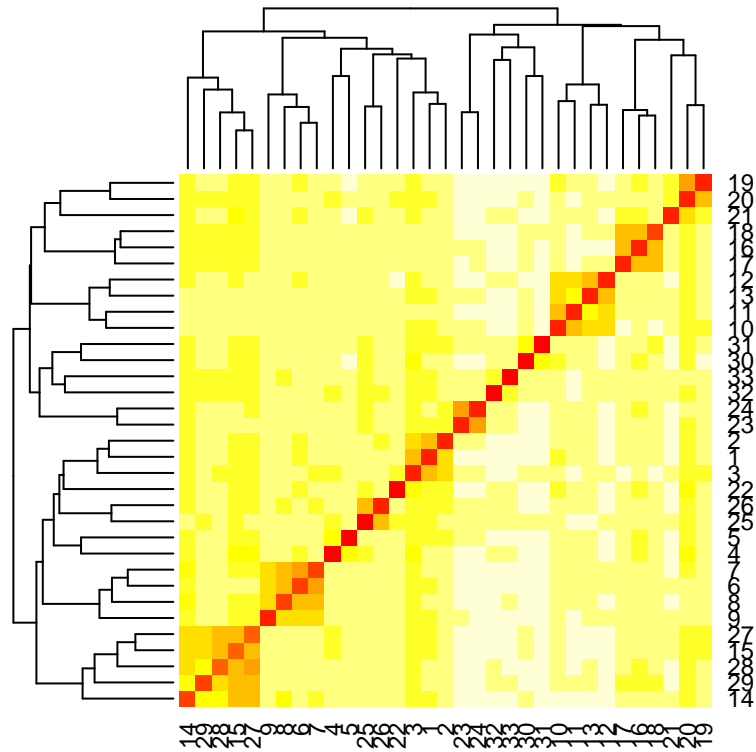
We use the `msa` library to align the sequences, then calculate the distance and draw some heatmaps as below:

```
#2.3
library(msa)

real_align <- msaClustalW("lizard_seqs.fasta",type="dna")

## use default substitution matrix

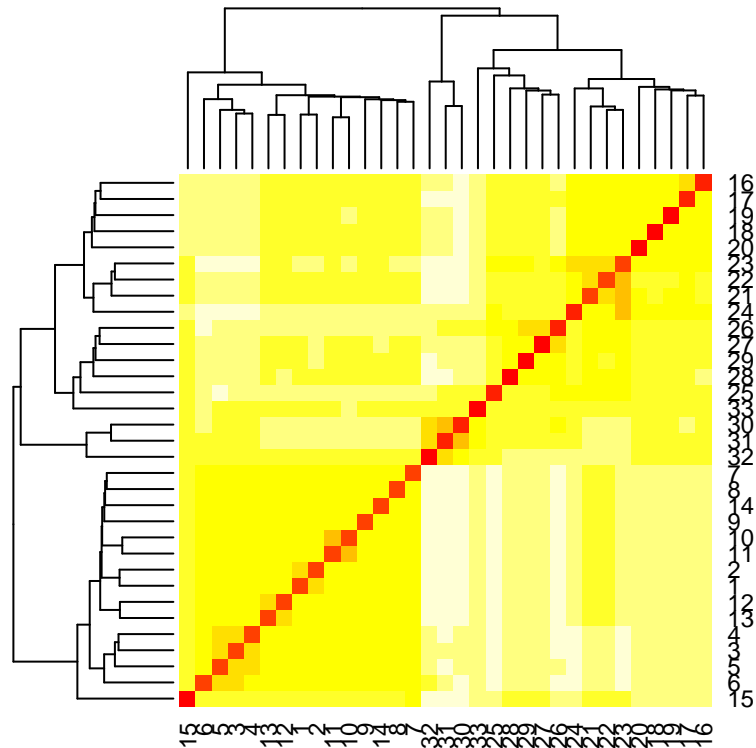
real_alignseq<- msaConvert(real_align, type="seqinr:alignment")
dist_real <- as.matrix(dist.alignment(real_alignseq, "identity"))
heatmap(dist_real)
```



```
ai1.1_align <- msaClustalW("AI_gene.fasta",type="dna")
```

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

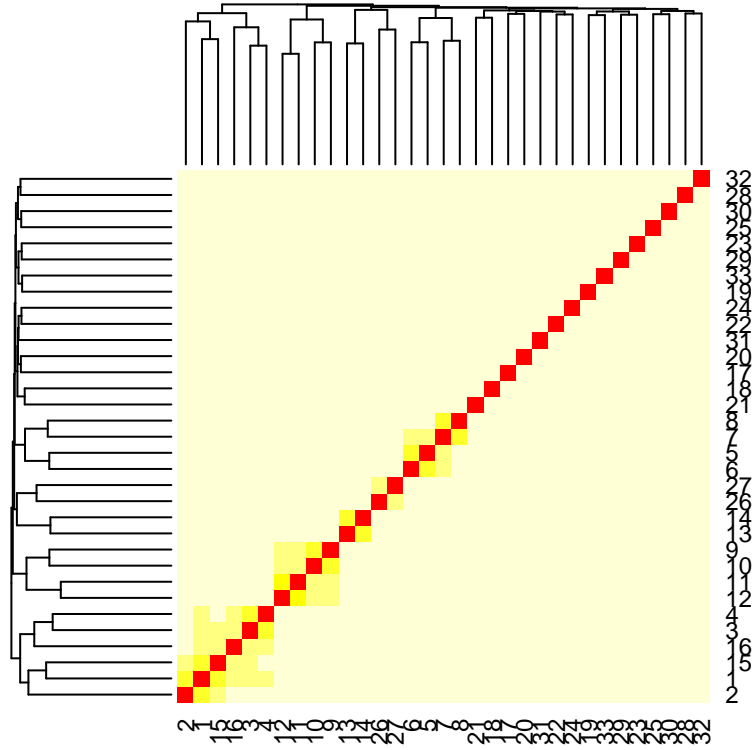
```
ai1.1_alignseq <- msaConvert(ai1.1_align, type="seqinr::alignment")
dist_a1.1 <- as.matrix(dist.alignment(ai1.1_alignseq, "identity"))
heatmap(dist_a1.1)
```



```
ai1.2_align <- msaClustalW("AI_gene2.fasta",type="dna")
```

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

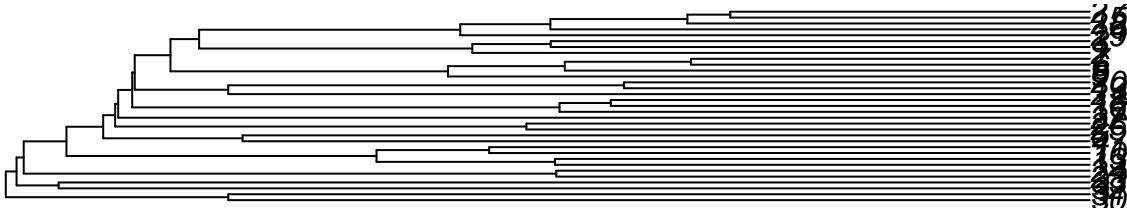
```
ai1.2_alignseq<- msaConvert(ai1.2_align, type="seqinr::alignment")
dist_a1.2 <- as.matrix(dist.alignment(ai1.2_alignseq, "identity"))
heatmap(dist_a1.2)
```



As we can see that the values in heatmap of AI sequences are quite low. It means that the AI gene has low (or no) connections with each other (because it randomly created) while the original gene is highly connected.

Question 3: Phylogeny reconstruction

UPGMA



NJ

