

732A51 Bioinformatics

# LAB 2 Bioinformatics

Hugo Morvan  
William Wiik

STIMA  
Department of Computer and Information Science  
Linköpings universitet

2024-11-22

Contents

<b>1</b>	<b>Question 1</b>	<b>1</b>
1.1	Question 1.1 . . . . .	2
1.2	Question 1.2 * . . . . .	3
<b>2</b>	<b>Question 2</b>	<b>4</b>
2.1	Question 2.1 . . . . .	4
2.2	Question 2.2* . . . . .	6
2.3	Question 2.3 . . . . .	6
<b>3</b>	<b>Question 3</b>	<b>6</b>
3.1	Question 3.1 . . . . .	6
3.2	Question 3.2* . . . . .	7

## 1 Question 1

In this exercise you will perform statistical analysis of three nucleotide data sets. First download the sequences from GenBank and save them in a fasta file. For this use the provided R script, 732A51\_BioinformaticsHT2023\_Lab02\_GenBankGetCode.R. This is a dataset of the RAG1 gene sequences from 33 lizard species. You are encouraged to read in detail the references in the script as they indicate many useful tools. Explore the dataset using the tools provided by the `ape` and `seqinr` packages. Take note of the lengths of all the sequences and the base composition.

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

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

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

## 1.1 Question 1.1

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

**Answer:**

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

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

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

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

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

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

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

```
## 33 DNA sequences in binary format stored in a list.
##
## Mean sequence length: 1982.879
##   Shortest sequence: 931
##   Longest sequence: 2920
##
## Labels:
## artificial_JF806202
## artificial_HM161150
## artificial_FJ356743
## artificial_JF806205
## artificial_JQ073190
## artificial_GU457971
## ...
##
## Base composition:
##   a       c       g       t
## 0.312 0.203 0.230 0.255
## (Total: 65.44 kb)
```

a: 31.2% c: 20.3% g: 23.0% t: 25.5%

We obtain almost the same base composition as the true base composition for lizard sequences.

## 1.2 Question 1.2 \*

## 2 Question 2

### 2.1 Question 2.1

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

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

calc_seq_stats <- function(all_seq, index_seq) {

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

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

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

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

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

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

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

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

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

}
```

```

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

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

# Simulate lizard sequence 1 (artifical_JF806202)
calc_seq_stats(all_seq = sim_sequences, index_seq = 1)

## [1] "Base composition for artifical_JF806202: a=0.286, c=0.207, g=0.257, t=0.251"
## [1] "GC content for artifical_JF806202: 0.464"
## [1] "CpG content for artifical_JF806202: 0.045"
## [1] "ApT content for artifical_JF806202: 0.07"

# =====
# From chatGPT, how to get frequency every charachter per sequence for an out file
# =====

process_protein_file <- function(file_path) {
  # Read the file
  lines <- readLines(file_path)

  # Extract headers (lines starting with ">")
  headers <- grep("^>", lines, value = TRUE)

  # Extract sequence lines (lines not starting with ">")
  header_indices <- grep("^>", lines)

  # Create a list of sequences
  sequences <- mapply(function(start, end) {
    paste(lines[(start + 1):(end - 1)], collapse = "")
  }, header_indices, c(header_indices[-1] - 1, length(lines)), SIMPLIFY = TRUE)

  # Remove ">" from headers
  headers <- gsub("^>", "", headers)

  # Create the data frame
  protein_data <- data.frame(
    Header = headers,
    Sequence = sequences,
    stringsAsFactors = FALSE
  )
}

```

```
# https://www.ebi.ac.uk/jdispatcher/st/emboss\_transeq/summary?jobId=emboss\_transeq-I20241120-223345-0023
transeq_true <- process_protein_file("transeq_true.out")
chars <- strsplit(transeq_true[1,2], NULL)[[1]] # JF806202
char_table <- prop.table(table(chars))
round(char_table[order(table(chars))],2)
```

```
## chars
##   X   Y   D   H   F   G   E   N   C   I   W   P   A   V   *   R
## 0.00 0.01 0.02 0.02 0.02 0.02 0.03 0.03 0.03 0.04 0.04 0.05 0.05 0.05 0.05 0.05
##   T   Q   K   M   S   L
## 0.05 0.06 0.06 0.07 0.09 0.15
```

```
# https://www.ebi.ac.uk/jdispatcher/st/emboss\_transeq/summary?jobId=emboss\_transeq-I20241120-221152-0336
transeq_sim <- process_protein_file("transeq_sim.out")
chars <- strsplit(transeq_sim[1,2], NULL)[[1]] # seq_1
char_table <- prop.table(table(chars))
round(char_table[order(table(chars))],2)
```

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

## 2.2 Question 2.2\*

## 2.3 Question 2.3

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

# 3 Question 3

## 3.1 Question 3.1

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



### 3.2 Question 3.2\*