# Exercises 1-3
# Statistics in Genetics

Alfons Edbom Devall

February 7, 2022

Statistics in Genetics, 7.5hp

# Contents

# 1 Chapter 1

## 1.1 Exercise 1

### 1.1.1 Background

A common first step when analyzing the genomic content of an organisms DNA is to check the percentage of the G- and C bases in the sequence of nucleotides, henceforth called the GC-content. If a region of the sequence has a higher GC-content than the background GC-content, it is likely that a gene/genes can be found there. Since calculating the GC-content is relatively fast and easy, with the potential to give good hints of interesting regions of the genome, this is often the first step when analyzing a new genome.
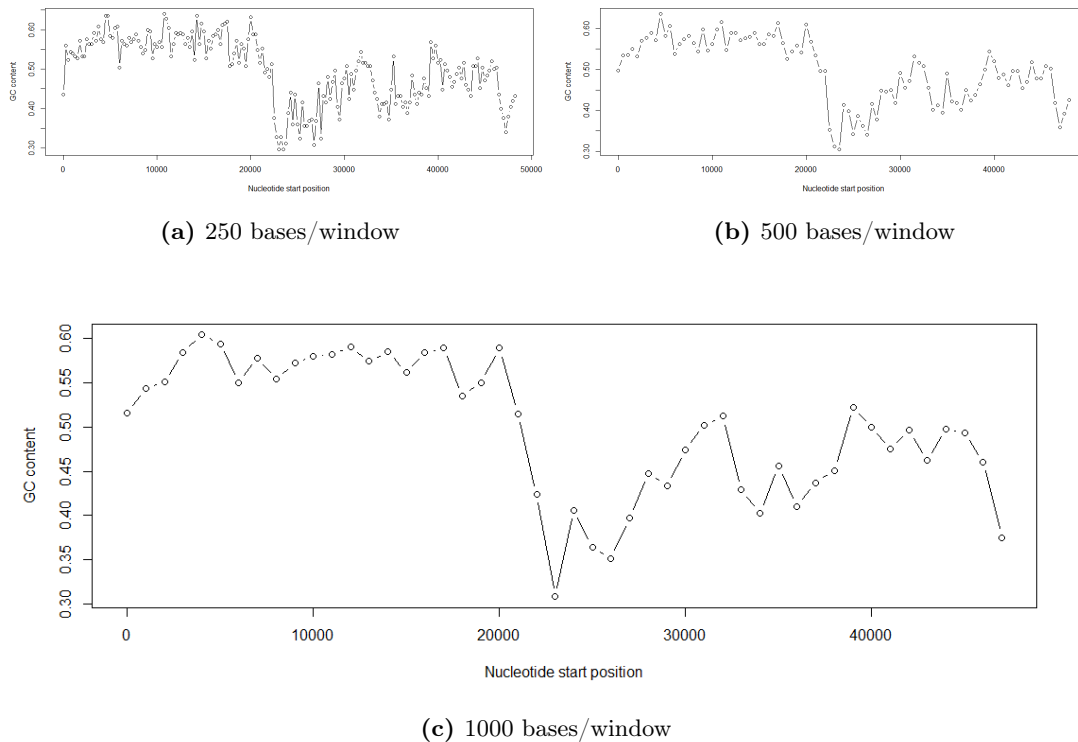
In this exercise the complete genomic sequence of the Enterobacteria phage lambda (NC_001417) was analyzed for its GC-content with various choices of window sizes.

### 1.1.2 What was done and discussion

In this exercise, the complete genome of Enterobacteria phage lamda was downloaded from GenBank. The genome was then analyzed by using the programming language R. First the FASTA-file was loaded using the seqinr package. Then a function to create a sliding window plot was made, called "slidingwindowplot", which takes as input both the window size to use and what sequence the plot should be made of. The effect of using different sizes of the window size then investigated.
To view the code for this exercise, see Appendix A.A

In figure 1 it can be seen that different sizes of the sliding window does appear to change which regions of high GC-content that are found in the genome. With a small size of the sliding window as in 1a, most regions with slightly higher GC-content than the background will be captured, but will instead be a bit cluttered and if they window size gets too small it starts to get too noisy. If using a larger size of the sliding window, as in figure 1c the data is getting more smooth, and as a consequence some of the smaller regions of high GC-content will be missed. To balance out these opposing factors, a window size somewhere in-between can be used, which then will give a good and high resolution figure, that is not too noisy and not too smooth, which can be seen in figure 1b.

**(a)** 250 bases/window



**(b)** 500 bases/window



**(c)** 1000 bases/window

**Figure 1** – Sliding window plots of the genome of Enterobacteria phage lamda using different sizes of the sliding window.

# 2   Chapter 2

## 2.1   Exercise 1

### 2.1.1   Background

When studying the genome of an organism, it is often of interest to identify regions of the DNA that codes for a gene.

To identify genes, a good start is to find the open reading frames (ORF) in the subject-DNA. An ORF consists of a DNA-sequence that starts with a start codon (e.g. ATG) and ends with a stop codon (TAA, TAG, TGA), where both the start and stop codon are in the same reading frame. The length of the ORF is then calculated by counting the number of nucleotides from the start codon to the stop codon, for it to be an ORF it needs to be divisible by 3 (since a codon consists of 3 nucleotides).

In this exercise the human(NC_001807), chimp( NC_001643) and mouse (NC_005089) mitochondrial DNA (mtDNA) has been investigated. The genetic code for mitochondria in vertebrates is slightly different than the standard one. The differences that are of interest for this exercise is that it has different start and stop codons: In mtDNA the

start codons are: ATG and ATA, and the stop codons are: TAA, TAG, AGA and AGG. All ORFs in the mtDNA have been identified and candidate protein coding genes has been selected. Finally, it was calculated how large fraction of the whole mtDNA is covered by the candidate genes.
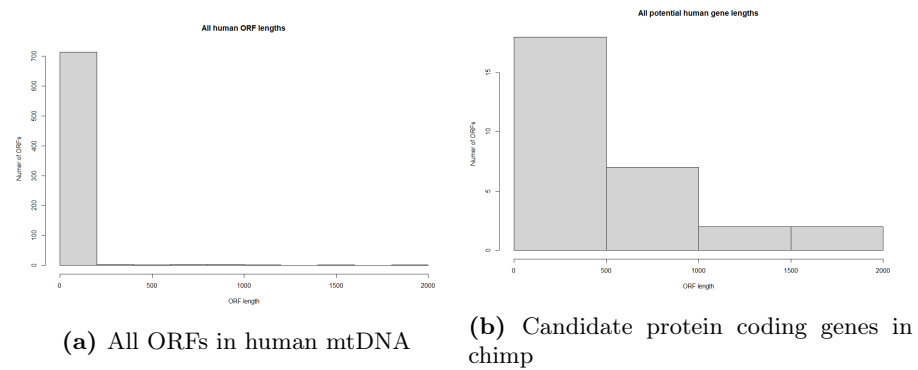
### 2.1.2   What was done and discussion

In this exercise the complete mtDNA from human, chimp and mouse was downloaded from GenBank. The mtDNA was then analyzed by using the programming language R. First the ORFik package was used to identify all ORFs in the sequence. A histogram of the lengths of the found ORFs was then done and can be seen in the figures 2a, 3a and 4a. Then the minimum length of the found ORFs was increased in small increments until I felt it was an appropriate number of ORFs left and also had an good approximation for how short a gene could potentially be. I chose 40 codons (120 base pairs) to be the minimum length to classify a found ORF as a candidate protein coding gene, and a histogram of the ORF lengths can be seen in the figures 2b, 3b and 4b.

When all potential protein coding genes had been identified, the fraction of the whole mtDNA sequence that was covered by the genes was calculated. This was done by summing up the lengths of all genes and then dividing by the whole sequence length. The resulting portions for human, chimp and mouse mtDNA was: 86.5%, 89.1% and 86.7% respectively.
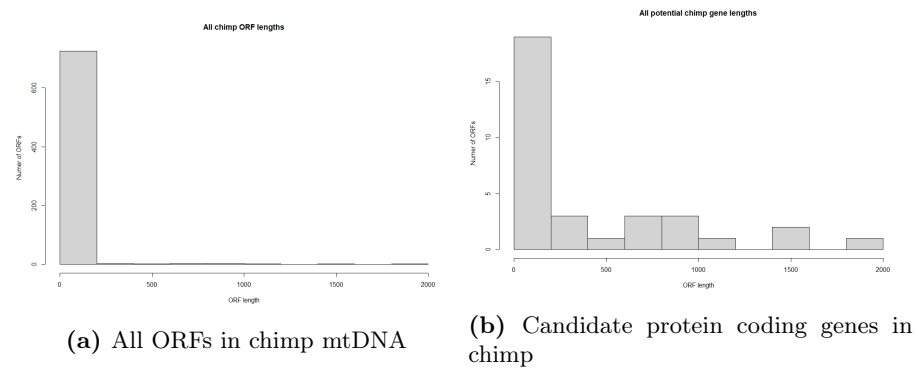
To solve this exercise a few different approaches could have been taken that may have given a more accurate result. When choosing what ORFs to be considered as genes, it might have been a better approach to generate random sequences of similar length to the mtDNA sequences and then seen how many ORFs at different lengths one could expect by chance and then choose a minimum length depending on that result.
When calculating what portion of the mtDNA that is covered by potential protein coding genes one could also take into account which regions of the actual mtDNA sequence is covered by each gene. This would most likely give a significantly more accurate reading of how large portion of the mtDNA sequence codes for genes. This is because the sequence is read in both directions(forwards and backwards) and thus contains some overlapping regions that right now is calculated twice.
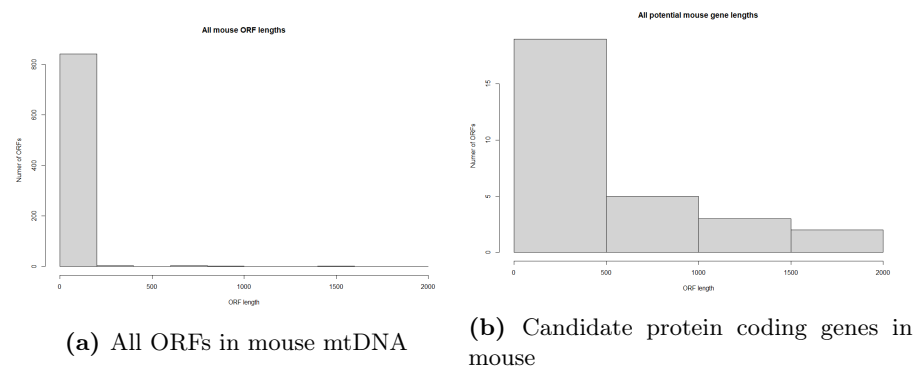To view the code for this exercise, see Appendix A.B

**(a)** All ORFs in human mtDNA

**(b)** Candidate protein coding genes in chimp

**Figure 2** – Histograms showing the ORFs found in human mtDNA as well as the lengths of the candidate protein coding genes



**(a)** All ORFs in chimp mtDNA

**(b)** Candidate protein coding genes in chimp

**Figure 3** – Histograms showing the ORFs found in chimp mtDNA as well as the lengths of the candidate protein coding genes.



**(a)** All ORFs in mouse mtDNA

**(b)** Candidate protein coding genes in mouse

**Figure 4** – Histograms showing the ORFs found in mouse mtDNA as well as the lengths of the candidate protein coding genes.

# 3 Chapter 3

## 3.1 Exercise 1

### 3.1.1 Background

It is often of interest to know how similar two DNA sequences are to each other. This similarity can then give hints about how closely related two sequences are and if they have a common ancestor. But with the large number of possible alignments, this task is not trivial, but with modern technology it is possible. There exists algorithms that can find both the best local alignment as well as the global alignment of two DNA sequences.

In this exercise the local and global alignments of the *D.melanogaster* gene *eyeless*(X79493) and the human gene *aniridia*(AY707088) has been investigated for their statistical significance.

### 3.1.2 What was done and discussion

In this exercise the eyeless gene and the aniridia gene was downloaded from GenBank. They were then analyzed by using the programming language R. First the R-package Biostrings was used to convert the downloaded FASTA files to strings. Then a local and a global alignment was done. The gap opening and gap extension penalty was 3 and 1 respectively, and the substitution matrix used for the alignment looks like this:

|   | A  | C  | G  | T  |
|---|----|----|----|----|
| A | 1  | -1 | -1 | -1 |
| C | -1 | 1  | -1 | -1 |
| G | -1 | -1 | 1  | -1 |
| T | -1 | -1 | -1 | 1  |

The resulting score for the local alignment was 183 for the score of the global alignment was -1204.
Next, the statistical significance was tested by creating 1000 random sequences of the same length as both the *eyeless* gene and the *aniridia*. Both a local and global alignment was done on all randomly generated sequences. The resulting scores can be seen in figure 5. Comparing the resulting scores from the random sequences and the score from the actual alignments we can be seen that the following: For the local alignment, the highest score obtained in the random sequences is 18, and the actual alignment score was 183, which is about a 10 times higher score, therefore the chance that this high of a score is very unlikely to be obtained just by chance. Similarly for the global alignment the highest score for a random sequence was -1442, while the score for the actual alignment was -1204, which is quite a bit higher than the random sequences. Therefore this score is also very unlikely to be this high or higher just by chance.
To view the code for this exercise, see Appendix A.C

**(a)** Global alignment of 1000 randomly generated sequences



**(b)** Local alignment of 1000 randomly generated sequences

**Figure 5**

A few different approaches could have been taken when solving this exercise. The largest improvement could potentially have been to use a more biologically-based transition matrix, like a BLOSUM or PAM matrix, as well as different values for the gap opening and gap extension penalties. I was however unsure of how far apart it is thought that these gene sequences are from each other, so I decided to just use a standard and basic transition matrix and gap penalties as to not use a matrix or values that captures some other biological factors that I was not interested in this exercise.

When generating the random sequences I chose to generate all of them to exactly the same length as both the *eyeless* and *aniridia* genes, and with a 25% chance for all nucleotides at each position. It may have given a more accurate result if the nucleotide percentages reflected more accurately the actual makeup of the genes, but I chose to continue using just a very simple model for this exercise.

# Appendices

## A    R-code

The sections below shows the R-code used to solve the exercises in this report. All code can also be found at this github-link: `https://github.com/AlfonsEdbom/Exercises_Statistics_in_Genetics.git`

### A.A    Exercise 1.1



```
Ex1.R ×    Joe_project_1.R ×    Joe_project_2.R ×    virus ×

     Source on Save

 1  require(pacman)  # Gives a confirmation message.
 2  pacman::p_load(pacman, seqinr, here)
 3
 4  #Get current path and get the path to the FASTA file
 5  here <- here("sequences", "NC_001416.fasta")
 6
 7  #Read the fasta file and get only the sequence
 8  virus <- read.fasta(here, forceDNAtolower = FALSE, set.attributes = FALSE) #Read the FASTA FILE
 9  virusseq <- virus[[1]]
10
11  #Function taken from https://a-little-book-of-r-for-bioinformatics.readthedocs.io/en/latest/src/chapter2_answers.html
12  #Creates a sliding windows plot of with the size windowsize
13  slidingwindowplot <- function(windowsize, inputseq)
14 ▾ {
15    starts <- seq(1, length(inputseq)-windowsize, by = windowsize)
16    n <- length(starts)
17    chunkGCs <- numeric(n)
18 ▾  for (i in 1:n) {
19      chunk <- inputseq[starts[i]:(starts[i]+windowsize-1)]
20      chunkGC <- GC(chunk)
21      chunkGCs[i] <- chunkGC
22 ▴  }
23    plot(starts,chunkGCs,
24        type="b",
25        xlab="Nucleotide start position",
26        ylab="GC content")
27 ▴ }
28
29  #Calls slidingwindowplot at specified windowsize and creates a plot
30  slidingwindowplot(250, virusseq)
31
32  #Close packages and plots
33  p_unload(all)
34  detach("package:datasets", unload = TRUE)
35
36  # Clear console
37  cat("\014")  # ctrl+L
38
```

**Figure 6** – The code used to create sliding window plots in Exercise 1.1

## A.B    Exercise 2.1



**Figure 7** – The code used to solve Exercise 2.1

## A.C   Exercise 3.1

```
1   require(pacman)  # Gives a confirmation message.
2   pacman::p_load(pacman, seqinr, here, Biostrings, universalmotif)
3   library(universalmotif)
4
5   #Get the file paths to sequences
6   fly_file <- here("sequences", "X79493.fasta")
7   human_file <- here("sequences", "AY707088.fasta")
8
9   #Load fasta files into DNA-strings
10  fly_sequence <- readDNAStringSet(fly_file)
11  human_sequence <- readDNAStringSet(human_file)
12
13  #Define substitution matrix
14  sub_mat <- nucleotideSubstitutionMatrix(match=1, mismatch = -1, baseOnly = TRUE)
15
16  #Make a local alignment
17  local_alignment <- pairwiseAlignment(pattern = fly_sequence,
18                      subject = human_sequence,
19                      type = "local",
20                      substitutionMatrix = sub_mat,
21                      gapOpening = 3,
22                      gapExtension = 1)
23
24  #Make a global alignment
25  global_alignment <- pairwiseAlignment(pattern = fly_sequence,
26                                  subject = human_sequence,
27                                  type = "global",
28                                  substitutionMatrix = sub_mat,
29                                  gapOpening = 3,
30                                  gapExtension = 1)
31
32
33  #Generate 1000 randomized of same length as human and fly sequence
34  fly_random <- create_sequences(alphabet = "DNA", seqnum = 1000, seqlen = width(fly_sequence))
35  human_random <- create_sequences(alphabet = "DNA", seqnum = 1000, seqlen = width(human_sequence))
36
37  #Saves the score of the local and global alignments
38  local_random <- numeric(length(human_random)) #Saves the score of the local alignment
39  global_random <- numeric(length(human_random))
40
41  for (i in 1:length(human_random)){
42      local_random[i] = pairwiseAlignment(pattern = fly_random[i],
43                                  subject = human_random[i],
44                                  type = "local",
45                                  substitutionMatrix = sub_mat,
46                                  gapOpening = 3,
47                                  gapExtension = 1,
48                                  scoreOnly = TRUE)
49
50      global_random[i] = pairwiseAlignment(pattern = fly_random[i],
51                                  subject = human_random[i],
52                                  type = "global",
53                                  substitutionMatrix = sub_mat,-
54                                  gapOpening = 3,
55                                  gapExtension = 1,
56                                  scoreOnly = TRUE)
57  }
58
59  #Creates histograms of the scores of
60  hist(local_random,
61      main = "Local alignment scores of random sequences",
62      xlab = "Score",
63      ylab = "Number")
64  hist(global_random,
65      main = "Global alignment scores of random sequences",
66      xlab = "Score",
67      ylab = "Number")
68
69  #Checks the probability of getting higher score than the global/local alignment
70  local_num_sign <- local_random[local_random >= score(local_alignment)]
71  local_pvalue <- local_num_sign/length(local_random)
72
73
74  global_num_sign <- global_random[global_random >= score(global_alignment)]
75  global_pvalue <- global_num_sign/length(global_random)
76
```

**Figure 8** – The code used to solve Exercise 3.1