

Genome Assembly

officialprofile

2021-09-05

Contents

1	Preface	5
1.1	Prerequisites	6
2	Introduction to genome assembly	7
3	Basic string manipulations	9
3.1	Generating random DNA sequence	9
3.2	Finding cDNA	10
3.3	Finding reverse complementary	11
4	Graph theory	13
5	De Bruijn graph	15
6	Burrows-Wheeler transform	17
6.1	Introduction	17
6.2	Burrows-Wheeler matrix	17
6.3	Inverse transform	19

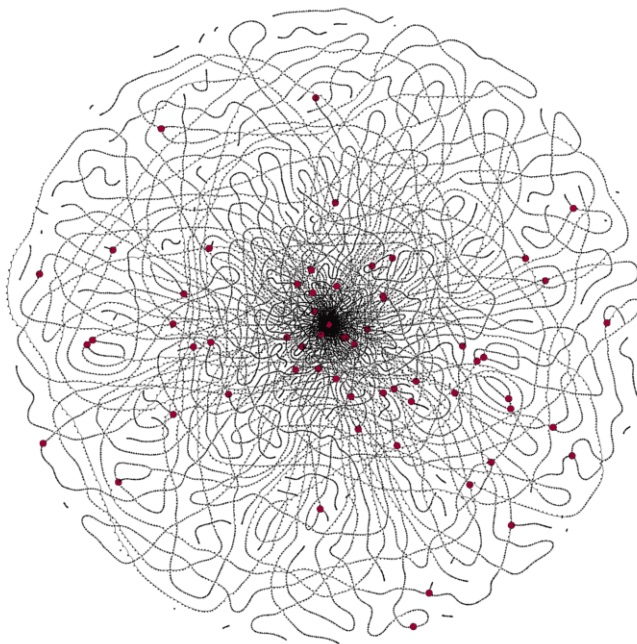
Chapter 1

Preface

This mini textbook describes selected algorithms that play a main role in the *de novo* genome assembly. The premise of this book is to construct these algorithms from the very bottom and explain step by step the main ideas that stand behind them.

We will try to avoid as much as possible ready-to-use implementations, which are, of course, available, but their use wouldn't serve the educational purpose. Naturally, many applications are included as well.

By default the code is written in R, but at certain points python is also mentioned.



A HiFi De Bruijn graph for a pile of reads from *Drosophila* genome sequencing. Each dot represents a k-mer ($k=23$), the edges denote neighboring k-mers. The larger red dots mark the head of heterozygous bubbles. Source: pacb.com.

1.1 Prerequisites

It is assumed that the reader

1. has a basic understanding of genetics;
2. has some experience with programming in R (is familiar with pipe syntax, etc.)
3. had some contact with higher mathematics, e.g. statistics, graph theory.

Throughout the book the following libraries are being used and it is assumed that the reader has them loaded.

```
library(stringr)
library(dplyr)
library(igraph)
```

Chapter 2

Introduction to genome assembly

Genome assembly has been regarded as one of the most important and most challenging problems in bioinformatics, at least since the late 80's, when the Human Genome Project was announced.

Genomes tend to be almost inconceivably long. Human DNA for example is approximately thousand times longer than the Bible, and some species have their genome even dozens times longer than ours. Yet the algorithmic part of the problem arises not mainly because of their length. This was definitely the case for the biotechnologists, who few decades ago were constrained within insufficiently effective sequencing techniques. From bioinformaticians perspective though the core of the problem is located deeper.

A good analogy, explaining the nature of the challenge, is a shredded newspaper. De novo assembly resembles reconstruction of the original document from a set of its unarranged pieces, like on the figure below.

As one can imagine the problem is difficult, but not because the newspaper has twenty or forty pages. If we had just one the assembly still wouldn't be easy. The length of the genome, although far from being irrelevant, is for us more of a secondary issue. On the other hand, repetitive patterns, which for biotechnologists are not a problem at all, will complicate our journey substantially. Reader will perhaps come up to the very similar conclusions as the succeeding topics will uncover.

One should also underline the fact that transcriptome assembly is not the same as the genome assembly. Read coverage of the latter is relatively uniform, whereas transcriptome can be differentially expressed and therefore the frequency of its reads can vary a lot. In genome assembly non-uniform abundance of the reads simply indicates existence of repetitions. In the case of transcription

product it is much less straightforward. Methods that overcome this issue exist but they are beyond of the scope of this textbook.

Chapter 3

Basic string manipulations

Before we dive into the main topics let's warm up with few basic exercises. On the one hand some of these problems might be regarded as a form of general familiarization with string manipulations. Being able to get through them is in a sense a must, at least from a perspective of algorithms that will be discussed later. On the other hand some of the problems, like generating k-mers or suffixes, are absolutely crucial in genome assembly and one has to know them to a tee. So, without further ado, let's do the warm up.

3.1 Generating random DNA sequence

Let's start with writing a function that, for a given integer $n > 0$, returns randomly generated DNA sequence.

```
random_sequence <- function(n){  
  return(sample(c('A','C','G','T'), n, replace = TRUE))  
}
```

```
random_sequence(10)  
#> [1] "G" "A" "C" "C" "G" "T" "G" "A" "T" "C"
```

If we want our returned sequence to be in a form of a single string we can use `paste()` function with `collapse` parameter equal to `' '`.

```
paste(random_sequence(10), collapse = ' ')  
#> [1] "TCAACACTCG"
```

For now we will try to stick to the first form, mainly because base R is not that good at manipulating strings (at this point we try to not employ any additional libraries like `stringr`). If a DNA is equal to `'ACGTG'` then we won't get a substring

'CG' simply by writing `DNA[2:3]` (one could do this in python). On the other hand, if `DNA` is a vector `c('A', 'C', 'G', 'T', 'G')` then `DNA[2:3]` will return `c('C', 'G')` that in turn can be easily merged into 'CG'.

3.2 Finding cDNA

One must know that C (cytosine) is complementary to G (guanine), and A (adenine) is complementary to T (thymine). In order to create a cDNA we will at first construct a list (in Python it would be a dictionary) of complementary nucleobases.

```
c_bases <- list('A' = 'T', 'C' = 'G', 'G' = 'C', 'T' = 'A')

c_bases['A'][[1]]
#> [1] "T"
```

This data structure will allow us to directly pull out complementary bases. Let's build a function that utilizes this.

```
complementary_sequence <- function(sequence){
  complementary = c()
  for (base in sequence){
    complementary = c(complementary, c_bases[base] [[1]])
  }
  return(complementary)
}
```

Please note that we append complementary nucleobase in a seemingly non-optimal way, namely `sequence = c(sequence, new_nucleobase)`. One could argue that this should be done with `append()`, i. e. `sequence = append(sequence, new_nucleobase)`, but in fact the `append` in R is regarded as a relatively slow function and our solution is usually more recommended.

```
seq <- random_sequence(20)
cat(' DNA =', seq, '\n')
#> DNA = G T C G C G C T C T A G A G G T C G G G
cseq <- complementary_sequence(seq)
cat(' cDNA =', cseq)
#> cDNA = C A G C G C G A G A T C T C C A G C C C
```

There is also a neater way to compute, and in a way avoid, these kind of loops. One can achieve this through functions like `apply()`, `sapply()` or `lapply()`.

```
complementary_sequence <- function(sequence){
  sapply(seq, function(base) c_bases[base] [[1]])
}
```

The clear downside of using functions like the one above is loss of legibility. Applying them can also be not as straightforward and intuitive as writing a conventional loop. Advanced programmers though find these functions very handy, fast and as readable as standard for or while loop. For that reason we will try to balance these two approaches out. We don't want our code to be too hermetic, but also we should not limit ourselves by deliberately avoiding legitimate solutions.

3.3 Finding reverse complementary

Obtaining reverse complementary is very similar. The only difference is that instead of appending complementary nucleobases we will prepend them.

```
reverse_complementary <- function(sequence){  
  reverse = c()  
  for (base in sequence){  
    reverse = c(c_bases[base][[1]], reverse)  
  }  
  return(reverse)  
}  
  
seq <- random_sequence(20)  
cat(' DNA =', seq, '\n')  
#> DNA = A G T C C A G C C G T G C A C T A C T T  
rcseq <- reverse_complementary(seq)  
cat(' cDNA =', rcseq)  
#> cDNA = A A G T A G T G C A C G G C T G G A C T
```


Chapter 4

Graph theory

Chapter 5

De Bruijn graph

Chapter 6

Burrows-Wheeler transform

The Burrows-Wheeler transform is one of the most effective lossless text compression method available. It provides a reversible transformation for text that makes it easier to compress. Of course, one may wonder what text compression has to do with genome assembly. As a matter of fact these two issues are closely related. But we should be more precise here - text compression is closely related to pattern matching which in turn is crucial for the genome assembly. In a broad sense compression algorithms look for patterns and try to remove repetitions. We want to take advantage of this feature, especially because repetitive patterns tend to be very abundant in genomic sequences.

It is worth mentioning that the Burrows-Wheeler transform is also closely related to suffix trees and suffix arrays, which are commonly used within pattern matching. This relationship will be studied later but perhaps reader should already keep the trivia in mind. (Donald Adjero, 2008)

6.1 Introduction

The Burrows-Wheeler transform method is often referred to as “block sorting”, because it takes a block of text and permutes it. By permuting a block of text we mean rearranging the order of its symbols. Once again, we should be more precise here because Burrows-Wheeler transform performs a specific type of permutation, namely *circular shift permutation*: all of the characters are moved one position to the left, and first character moves to the last position.

6.2 Burrows-Wheeler matrix

Consider the following sequence:

```
1 sequence <- 'GATTACA'
```

In order to create the Burrows-Wheeler matrix, from which the transform itself can be obtained, for the given string we at first add the dollar sign \$ at the end of the sequence.

```
1 sequence <- str_c(sequence, '$')
```

Afterwards we perform a series of circular shift permutations.

```
1 sequences <- c(sequence)
2 n          <- nchar(sequence)
3
4 for (i in 1:(n-1)){
5   sequence <- str_c(str_sub(sequence, 2, n),
6                     str_sub(sequence, 1, 1))
7
8   sequences <- c(sequences, sequence)
9 }
10
11 cat(sequences, sep = '\n')
12 #> GATTACA$
13 #> ATTACA$G
14 #> TTACA$GA
15 #> TACA$GAT
16 #> ACA$GATT
17 #> CA$GATTA
18 #> A$GATTAC
19 #> $GATTACA
```

Then we sort these sequences with the assumption that the dollar sign precedes lexicographically every other symbol.

```
1 sequences <- sort(sequences)
2 cat(sequences, sep = '\n')
3 #> $GATTACA
4 #> A$GATTAC
5 #> ACA$GATT
6 #> ATTACA$G
7 #> CA$GATTA
8 #> GATTACA$
9 #> TACA$GAT
10 #> TTACA$GA
```

For our convenience let's split these permutations into vectors of single characters.

```

1 bw.matrix      <- data.frame(matrix(, n, n))
2 colnames(bw.matrix) <- 1:n
3
4 for (i in 1:n){
5   bw.matrix[i, ] <- strsplit(sequences[i], split = '')[[1]]
6 }
7
8 knitr::kable(bw.matrix)

```

1	2	3	4	5	6	7	8
\$	G	A	T	T	A	C	A
A	\$	G	A	T	T	A	C
A	C	A	\$	G	A	T	T
A	T	T	A	C	A	\$	G
C	A	\$	G	A	T	T	A
G	A	T	T	A	C	A	\$
T	A	C	A	\$	G	A	T
T	T	A	C	A	\$	G	A

Thus we have created the **Burrows-Wheeler matrix**. Sequence in the last column is called the **Burrows-Wheeler transform**.

```

1 transform <- paste(bw.matrix[,n], collapse = '')
2
3 cat('The Burrows-Wheeler transform of',
4     sequence, 'is', transform)
5 #> The Burrows-Wheeler transform of $GATTACA is ACTGA$TA

```

6.3 Inverse transform

As we said at the very beginning the transform is reversible. Having only the transformed sequence we are going to reconstruct the Burrows-Wheeler matrix and initial sequence itself.

Firstly let's sort the characters of the transformed sequence.

```

1 first.sequence <- strsplit(transform, split = '')[[1]] %>% sort
2 paste(first.sequence, collapse = '')
3 #> [1] "$AAACGTT"

```

Note that this string is equivalent to the first column of the Burrows-Wheeler transform.

```

1 bw.inverse      <- data.frame(matrix(, n, 2))
2 colnames(bw.inverse) <- c(n, 1)
3

```

8	1
A	\$
C	A
T	A
G	A
A	C
\$	G
T	T
A	T

```
1 kmers <- apply(bw.inverse, 1,
2               function(x) paste(x, collapse = ''))
3 kmers
4 #> [1] "A$" "CA" "TA" "GA" "AC" "$G" "TT" "AT"
```

```
1 kmers <- sort(kmers)
2 kmers
3 #> [1] "$G" "A$" "AC" "AT" "CA" "GA" "TA" "TT"
```

```
1 supply(kmers, function(x) str_sub(x, 2, 2),
2       simplify = TRUE, USE.NAMES = FALSE)
3 #> [1] "G" "$" "C" "T" "A" "A" "A" "T"
```

[illegible]

```

7   colnames(bw.inverse)[i+1] = i
8   }
9   knitr::kable(bw.inverse)

```

8	1	2	3	4	5	6	7
A	\$	G	A	T	T	A	C
C	A	\$	G	A	T	T	A
T	A	C	A	\$	G	A	T
G	A	T	T	A	C	A	\$
A	C	A	\$	G	A	T	T
\$	G	A	T	T	A	C	A
T	T	A	C	A	\$	G	A
A	T	T	A	C	A	\$	G

Finally we move first column to the very end

```

1   bw.inverse[,n+1] <- bw.inverse[, 1]
2   bw.inverse      <- bw.inverse[,2:(n+1)]
3   colnames(bw.inverse)[n] = n
4
5   knitr::kable(bw.inverse)

```

1	2	3	4	5	6	7	8
\$	G	A	T	T	A	C	A
A	\$	G	A	T	T	A	C
A	C	A	\$	G	A	T	T
A	T	T	A	C	A	\$	G
C	A	\$	G	A	T	T	A
G	A	T	T	A	C	A	\$
T	A	C	A	\$	G	A	T
T	T	A	C	A	\$	G	A

One can also verify that `bw.matrix` and `bw.inverse` are in fact the same.

```

1   knitr::kable(bw.inverse == bw.matrix)

```

1	2	3	4	5	6	7	8
TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE
TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE
TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE
TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE
TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE
TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE
TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE
TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE

Additionally we can encapsulate the Burrows-Wheeler transform in a form of a single function.

```

1 BWT <- function(sequence){
2   sequence <- str_c(sequence, '$')
3   sequences <- c(sequence)
4   n <- nchar(sequence)
5
6   for (i in 1:(n-1)){
7     sequence <- str_c(str_sub(sequence, 2, n),
8                       str_sub(sequence, 1, 1))
9     sequences <- c(sequences, sequence)
10  }
11  sequences <- sort(sequences)
12
13  bw.matrix <- data.frame(matrix(, n, n))
14  colnames(bw.matrix) <- 1:n
15
16  for (i in 1:n){
17    bw.matrix[i, ] <- strsplit(sequences[i], split = '')[[1]]
18  }
19  return(paste(bw.matrix[,n], collapse = ''))
20 }

```

```

1 BWT('GATTACA')
2 #> [1] "ACTGA$TA"

```

One can verify that this output is equal to result we obtained earlier.

Out of pure curiosity lets check the Burrows-Wheeler transform for a longer sequence.

```

1 BWT('ATGCTCGTGCCATCATATAGCGGCGCGGATCTCTACGCGCG')
2 #> [1] "GTTTCCG$TCGGGGGAGGGTTGTCTCCCCCATCCAAACCAGA"

```

Please note that the input string has no identical characters at adjacent positions, whereas in the transformed sequence such situation appears quite often. These substrings of identical characters will allow us represent the sequence in a more condensed manner and expediate pattern matching.

Bibliography

Donald Adjero, Tim Bell, A. M. (2008). *The Burrows-Wheeler Transform: Data Compression, Suffix Arrays, and Pattern Matching*. Springer.