

Solutions Exercise Set 3

Problem 1: Suppose we have following fragments:

f1 = ATCCGTTGAAGCCGCGGGC
f2 = TTAACTCGAGG
f3 = TTAAGTACTGCCCG
f4 = ATCTGTGTCGGG
f5 = CGACTCCC GACACA
f6 = CACAGATCCGTTGAAGCCGCGGG
f7 = CTCGAGTTAAGTA
f8 = CGCGGGCAGTACTT

We know that the total length of the target molecule is about 50bp and may be ready to accept a solution of length between 45 and 55 bp. Assemble these fragments and obtain a consensus sequence. Be prepared to deal with errors. You may also have to use the reverse complement of some of the fragments.

Solution:

Fragments	Reversed and complemented
f1 = ATCCGTTGAAGCCGCGGGC	GCCCCGGCTTCAACGGAT
f2 = TTAACTCGAGG	CCTCGAGTTAA
f3 = TTAAGTACTGCCCG	CGGGCAGTACTTAA
f4 = ATCTGTGTCGGG	CCCGACACAGAT
f5 = CGACTCCC GACACA	TGTGTCGGGAGTCG
f6 = CACAGATCCGTTGAAGCCGCGGG	CCCGCGGCTTCAACGGATCTGTG
f7 = CTCGAGTTAAGTA	TACTTAACTCGAG
f8 = CGCGGGCAGTACTT	AAGTACTGCCCGCG

Starting with f2' and f3

CCTCGAGTTAA

TTAAGTACTGCCCG

CCTCGAGTTAAGTACT GCCCG

AAGTACTGCCCGCG

(f8')

CCTCGAGTTAAGTACTGCCCGCG

CTCGAGTTAAGTACT

(f7)

CCTCGAGTTAACTACTGCCCGCG

GCCCCGCGGCTTCAACGGAT

(f1')

CCTCGAGTTAACTACTGCCCGCG

CCCGCGGGCTTCAACGGAT

(f6')

CCTCGAGTTAACTACTGCCCGCG

CCTCGAGTTAACTACTGCCGCGGCTTACACGG	ATCTGTG	(last superstring)
	ATCTGTGT	CGGG (f4)
	TGTGT	CGGGAGTCG (f5)
<hr/>		
CCTCGAGTTAACTACTGCCGCGGCTTACACGGATCTGTGTGGAGTCG		

Problem 2: Let $F=\{\text{ATC, TCG, AACG}\}$. Find the best layout for this collection according to the sequence reconstruction model, with error level $e=0.1$. The same problem for $e=0.25$. Be sure to consider also reverse complements.

Solution:

With $e = 0.1$

- For $f_1 = \text{ATC}$ it matches well with substring $a = \text{ATC}$, $d(f_1, a) = 0 < 0.3$ (because $|f_1| = 3$ and $e^* |f_1| = 0.3$).
- For $f_2 = \text{TCG}$ it matches well with substring $a = \text{ATC}$, $d(f_2, a) = 0 < 0.3$ (because $|f_2| = 3$ and $e^* |f_2| = 0.3$).
- For $f_3 = \text{AACG}$, $e^* |f_3| = 0.4$.

If $f_3 = \text{AACG}$ is considered and it is matched with substring $a = \text{ATCG}$, then $d(f_3, a) = 1 \not< 0.4$ (because $|f_3| = 4$).

Taking reverse complement $f_3' = \text{CGTT}$.

f_3' matches well with $a = \text{CGTT}$ (see the table with error level 0.1), $d(f_3', a) = 0$.

<table border="1" style="border-collapse: collapse; margin-bottom: 10px;"> <tr><td>A</td><td>T</td><td>C</td><td></td></tr> <tr><td></td><td>T</td><td>C</td><td>G</td></tr> <tr><td>A</td><td>A</td><td>C</td><td>G</td></tr> <tr style="background-color: yellow;"><td>A</td><td>T</td><td>C</td><td>G</td></tr> </table> <p>Error level 0.25 One error is allowed.</p>	A	T	C			T	C	G	A	A	C	G	A	T	C	G	<table border="1" style="border-collapse: collapse; margin-bottom: 10px;"> <tr><td>A</td><td>T</td><td>C</td><td></td><td></td><td></td></tr> <tr><td></td><td>T</td><td>C</td><td>G</td><td></td><td></td></tr> <tr><td></td><td></td><td>C</td><td>G</td><td>T</td><td>T</td></tr> <tr><td style="background-color: yellow;">A</td><td style="background-color: yellow;">T</td><td style="background-color: yellow;">C</td><td style="background-color: yellow;">G</td><td style="background-color: yellow;">T</td><td style="background-color: yellow;">T</td></tr> </table> <p>Error level 0.1 No errors are allowed.</p>	A	T	C					T	C	G					C	G	T	T	A	T	C	G	T	T
A	T	C																																							
	T	C	G																																						
A	A	C	G																																						
A	T	C	G																																						
A	T	C																																							
	T	C	G																																						
		C	G	T	T																																				
A	T	C	G	T	T																																				

Considering $e = 0.25$.

- For $f_1 = \text{ATC}$ it matches well with substring $a = \text{ATC}$, $d(f_1, a) = 0 < 0.75$ (because $|f_1| = 3$ and $e^* |f_1| = 0.75$).
 - For $f_2 = \text{TCG}$ it matches well with substring $a = \text{ATC}$, $d(f_2, a) = 0 < 0.75$ (because $|f_2| = 3$ and $e^* |f_2| = 0.75$).
 - For $f_3 = \text{AACG}$,
- If $f_3 = \text{AACG}$ is considered and it is matched with substring $a = \text{ATCG}$, then $d(f_3, a) = 1 = 1$ (because $|f_3| = 4$ and $e = 0.25$, $e^* |f_3| = 1$).

Problem 3:

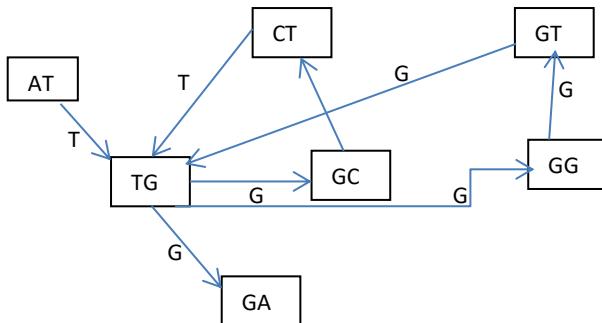
- a) You want to use the sequencing by hybridization method (SBH) to sequence a DNA fragment. For this, you are using a DNA array containing all DNA sequences of length 3 and test which of these sequences bind to your target. As a result, you find out that the target sequence has the following substrings of length 3:

{ ATG, CTG, GCT, GGT, GTG, TGA, TGC, TGG } Find at least 2 DNA sequences validating this data.

- b) How many solutions do you have if, using a DNA array containing all sequences of length 4, you obtain that the target sequence has the following substrings of length 4:
{ATGG, CTGA, GCTG, GGTG, GTGC, TGCT, TGGT } ?

Solution:

(a)



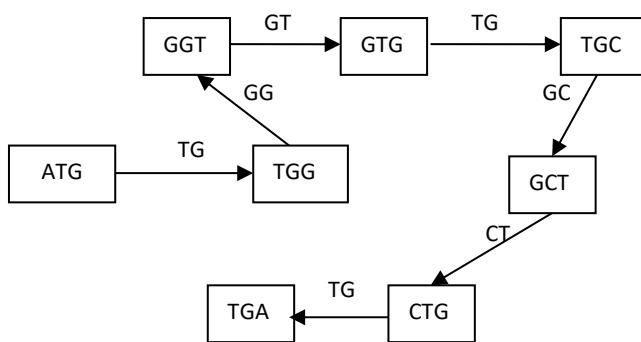
Following the Eulerian paths

$AT \rightarrow TG \rightarrow GC \rightarrow CT \rightarrow TG \rightarrow GG \rightarrow GT \rightarrow TG \rightarrow GA$ One possible sequence ATGCTGGTGA

$AT \rightarrow TG \rightarrow GG \rightarrow GT \rightarrow TG \rightarrow GC \rightarrow CT \rightarrow TG \rightarrow GA$ Another possible sequence ATGGTGCTGA

(b)

{ATGG, CTGA, GCTG, GGTG, GTGC, TGCT, TGGT}



One Eulerian path possible

$ATG \rightarrow TGG \rightarrow GGT \rightarrow GTG \rightarrow TGC \rightarrow GCT \rightarrow CTG \rightarrow TGA$

Solution:

ATGGTGCTGA

Problem 4: You are assembling a DNA sequence containing a repeat of the form XXX. Having given the fragments ATG, CTTGAT, TGT, TGTCA, TCAGAT, TGTAAC, find at least two such DNA sequences knowing that

- No fragment is included into some other
- The fragments provide “good linkage”, in the sense that all fragments (except the one covering the ends of the sequence) overlap with at least one fragment at left and with another at right.

Solution:

Fragments	Reversed and complemented
ATG	CAT
CTTGAT	ATCAAG
TGT	ACA
TGTCA	TGACA
TCAGAT	ATCTGA
TGTAAC	AGTTACA

One possible solution:

```

A T C A A G
- - - - A G T T A C A - - - - - - - - - - - -
- - - - - - - - C A T - - - - - - - - - - - -
- - - - - - - - - T G T - - - - - - - - - -
- - - - - - - - - - T G A C A - - - - - - - -
- - - - - - - - - - - - - - A T C T G A
A T C A A G T T A C A T G T G A C A T C T G A

```

Other possible solutions :

Fragment 4 and 5 – Alignment 1

T	G	T	C	A	-	-	-
—	—	T	C	A	G	A	T
T	G	T	C	A	G	A	T

Fragments 2 and 6 – Alignment 2

```

      - - - - - C T T G A T
      T G T A A C T - - -
      T G T A A C T T G A T

```

Consider above two alignments. It is clear that TGT must be part of the repeat.

One solution could be :

T G T C A G A T T G T A A A C T T T G A T - -
 - - - - - - - - - - - - - - - - A T G
 - - - - - - - - - - - - - - - T G T

Another solution :

More possibilities:

Result 3: ATG + Alignment 1 + TGT + Alignment 2: **ATGTCAGATGTGTAACTTGAT** – length 21.

Result 4: ATG + Alignment 2 + TGT + Alignment 1: **ATGTAACATTGATGTGTCAGAT** – length 21.

Problem 5.

- a) Assemble the following error-free fragments using the shotgun approach: ATGTG, GCCGCA, GTGCCG, TGTGCC.
b) The same problem as above, replacing the second fragment above with CCCGCA. Assemble the fragments using the shotgun approach. ~~Assemble the fragments using also the SBH-style shotgun approach. Compare the results and also with the result obtained at a), knowing that fragment CCCGCA had one substitution error – the correct one was the second fragment in a).~~

Solution:

(a)

$$\begin{array}{ccccccccc} & T & G & T & G & C & C & - & - \\ - & & G & T & G & C & C & & \\ A & T & G & T & G & & & & \\ - & - & - & - & G & C & C & G & C & A \\ \hline A & T & G & T & G & C & C & G & C & A \end{array}$$

(b)

Replacing the second fragment with CCCGCA.

$$\begin{array}{ccccccccc} & - & - & - & T & G & T & G & C & C \\ & - & - & - & G & T & G & C & C & G \\ & & & A & T & G & T & G & & \\ \hline C & C & C & G & C & A & - & - & - & C & G \\ \hline C & C & C & G & C & A & T & G & T & G & C & C & G \end{array}$$