

Problem (1.a)

The problem setup is just like on Homework 1, the only difference is you have an added assumption:

$$P(BB) = 0.$$

Thus, following from homework 1 / lecture:

$$P(E_1, \dots, E_N | AB) = \left(\frac{1}{2}\right)^N$$

$$\begin{aligned} P(E_1, \dots, E_N | AA) &= \prod_{i: \text{obs } A} P(E_i = 0 | AA) \prod_{j: \text{obs } B} P(E_j = 1 | AA) \\ &= P(E_i = 0 | AA)^{N_A} P(E_j = 1 | AA)^{N_B} \\ &= (1-e)^{N_A} e^{N-N_A} \end{aligned}$$

where $N_A = \# \text{ of } A\text{'s we observe}$, $N_B = \# \text{ of } B\text{'s}$.

So,

$$\begin{aligned} P(AB | E_1, \dots, E_N) &= \frac{P(E_1, \dots, E_N | AB) P(AB)}{P(E_1, \dots, E_N | AB) P(AB) + P(E_1, \dots, E_N | AA) P(AA)} \\ &= \frac{\left(\frac{1}{2}\right)^N P(AB)}{\left(\frac{1}{2}\right)^N P(AB) + (1-e)^{N_A} e^{N-N_A} P(AA)} \end{aligned}$$

Problem (1.b)

Same as (a), just a different numerator!

$$\begin{aligned} P(AA | E_1, \dots, E_N) &= \frac{P(E_1, \dots, E_N | AA) P(AA)}{P(E_1, \dots, E_N | AA) P(AA) + P(E_1, \dots, E_N | AB) P(AB)} \\ &= \frac{(1-e)^{N_A} e^{N-N_A} P(AA)}{\left(\frac{1}{2}\right)^N P(AB) + (1-e)^{N_A} e^{N-N_A} P(AA)} \end{aligned}$$

Problem 1c)

$$P(E_1, \dots, E_N | AA) = \prod_{i:A} P(E_i = 0 | AA) \prod_{i:B} P(E_i = 1 | AA)$$

$$= (1-e)^{N_A} e^{N-N_A} = \left(\frac{1-e}{e}\right)^{N_A} e^N$$

$$P(AA | E_1, \dots, E_N) = \frac{\left(\frac{1-e}{e}\right)^{N_A} e^N \left(\frac{1}{2}\right)}{\left(\frac{1-e}{e}\right)^{N_A} e^N \frac{1}{2} + \left(\frac{1}{2}\right)^N \frac{1}{2}}$$

$$= \frac{\left(\frac{1-e}{e}\right)^{N_A} e^N}{\left(\frac{1-e}{e}\right)^{N_A} e^N + \left(\frac{1}{2}\right)^N} = \alpha \Rightarrow \frac{\left(\frac{1-e}{e}\right)^{N_A} e^N}{\left(\frac{1-e}{e}\right)^{N_A} e^N} + \frac{\left(\frac{1}{2}\right)^N}{\left(\frac{1-e}{e}\right)^{N_A} e^N} = \frac{1}{\alpha}$$

$$\Rightarrow \frac{\frac{1}{2^N}}{\left(\frac{1-e}{e}\right)^{N_A} e^N} = \frac{1}{\alpha} - 1 \Rightarrow \frac{1}{\left(\frac{1-e}{e}\right)^{N_A} e^N 2^N} = \frac{1-\alpha}{\alpha}$$

$$\Rightarrow \frac{\alpha}{(1-\alpha) e^N 2^N} = \left(\frac{1-e}{e}\right)^{N_A} \Rightarrow \frac{\log\left(\frac{\alpha}{(1-\alpha) (2e)^N}\right)}{\log\left(\frac{1-e}{e}\right)} = N_A$$

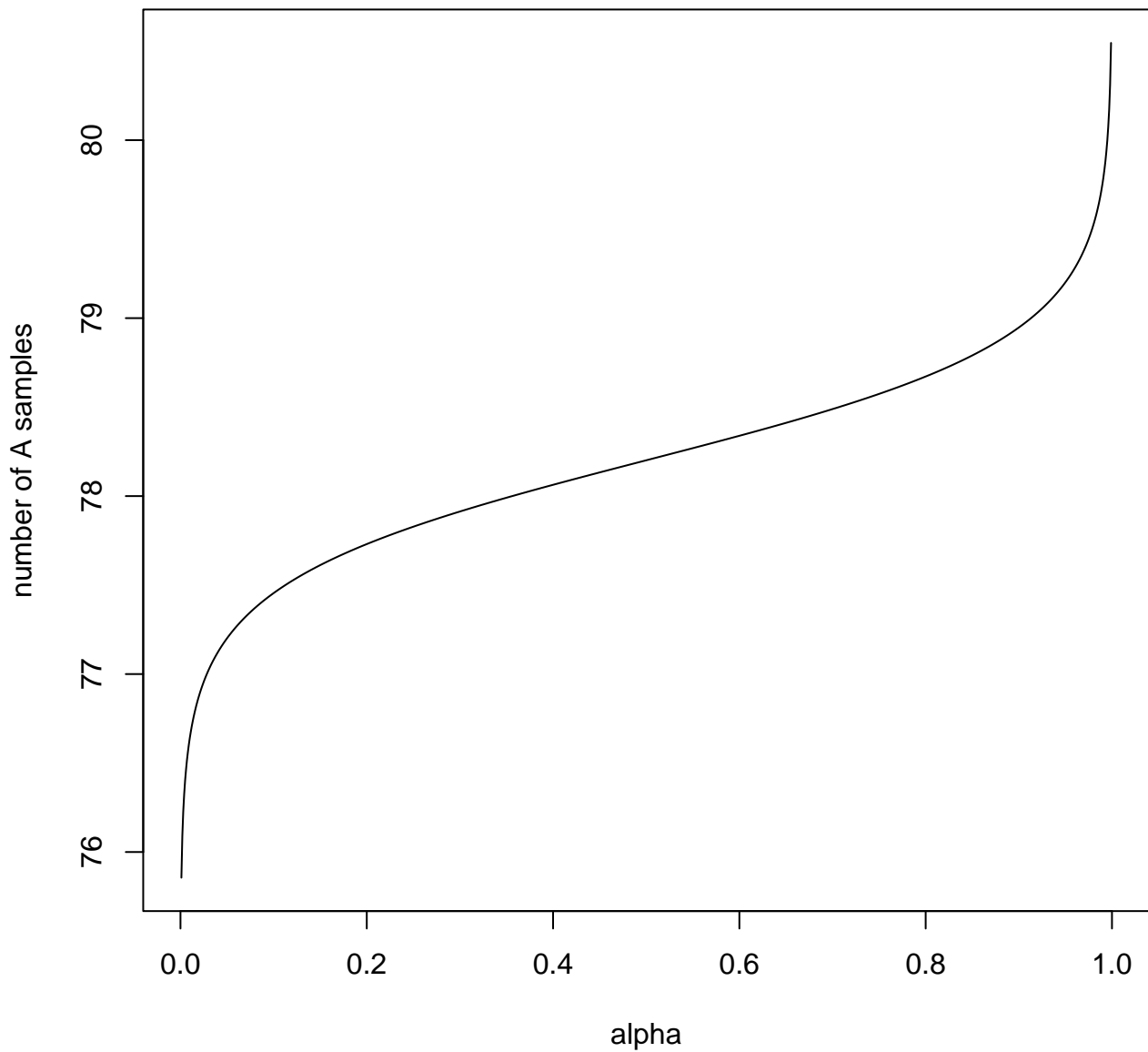
$$\Rightarrow N_A = \frac{\log(\alpha) - \log(1-\alpha) - N \log(2e)}{\log(1-e) - \log(e)}$$

Problem (1.2)

Our solution from (1c) is a function of α , N , and e :

$$N_A = \frac{\log(\alpha) - \log(1-\alpha) - N \log(ze)}{\log(1-e) - \log(e)}.$$

Since N and e are given, we can simply plot N_A as a function of α :



Problem (1e)

We need to account for case BB, which from class/homework is:

$$\begin{aligned}
 P(E_1, \dots, E_N | BB) &= \prod_{i \in A} P(E_i = 1 | BB) \prod_{j \in B} P(E_j = 0 | BB) \\
 &= (1-e)^{N-N_A} e^{N_A} = (1-e)^{N-N_A} e^{N_A}
 \end{aligned}$$

Note, given the description, we know $N_A = 0$, so,

$$P(E_1, \dots, E_N | BB) = (1-e)^N$$

$$\begin{aligned}
 P(BB | E_1, \dots, E_N) &= \frac{(1-e)^N P(BB)}{(1-e)^N P(BB) + \left(\frac{1}{2}\right)^N P(AB) + \left(\frac{1-e}{e}\right)^{N_A} e^N P(AA)} \\
 &= \frac{(1-e)^N P(BB)}{(1-e)^N P(BB) + \left(\frac{1}{2}\right)^N P(AB) + e^N P(AA)} = \alpha
 \end{aligned}$$

flip both sides:

$$1 + \left(\frac{\frac{1}{2}}{1-e}\right)^N \frac{P(AB)}{P(BB)} + \left(\frac{e}{1-e}\right)^N \frac{P(AA)}{P(BB)} = \frac{1}{\alpha}$$

$$\Rightarrow \left(\frac{1}{(1-e)2} \right)^N \frac{P(AB)}{P(BB)} + \left(\frac{e}{1-e} \right)^N \frac{P(AA)}{P(BB)} = \frac{1-\alpha}{\alpha}$$

Log both sides

$$\Rightarrow N \left[\log 1 - \log(2(1-e)) \right] + \log \frac{P(AB)}{P(BB)} + N \log \left(\frac{e}{1-e} \right)$$

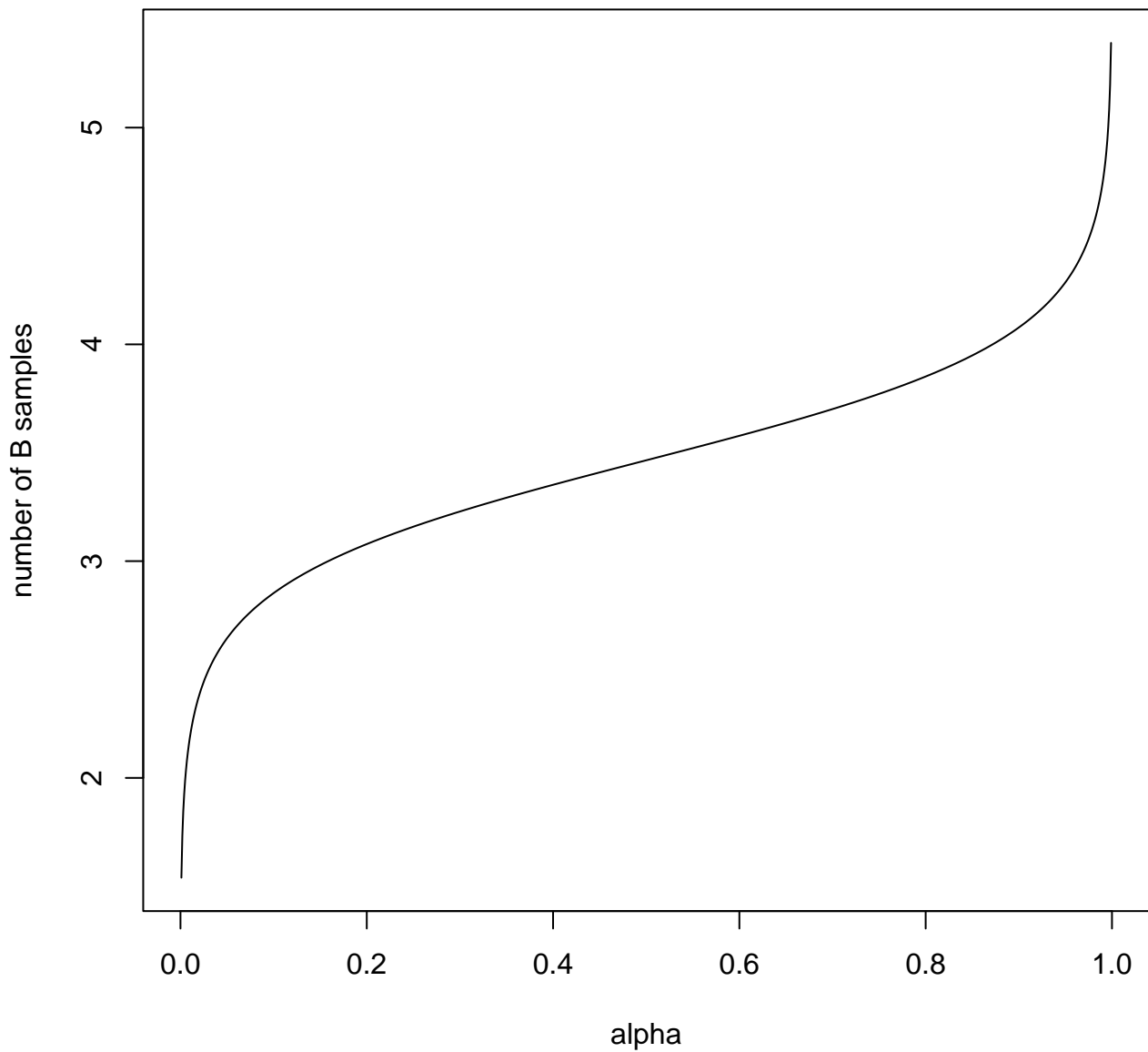
$$+ \log \frac{P(AA)}{P(BB)} = \log \left(\frac{(1-\alpha)}{\alpha} \right)$$

$$\Rightarrow N \left[\log \left(\frac{e}{1-e} \right) - \log(2(1-e)) \right] = \log \left(\frac{(1-\alpha)}{\alpha} \right) - \log \frac{P(AB)}{P(BB)} - \log \frac{P(AA)}{P(BB)}$$

$$\Rightarrow N = \frac{\log \left(\frac{(1-\alpha)}{\alpha} \right) - \log \left(\frac{P(AB)}{P(BB)} \right) - \log \left(\frac{P(AA)}{P(BB)} \right)}{\log \left(\frac{e}{1-e} \right) - \log(2(1-e))}$$

Fixing $e = 0.05$, $P(AA) = P(AB) = 0.4995$, $P(BB) = 0.001$

You see the plot:



Problem 2

Problem (2.a)

$t_1 = \text{AAAA CGCGAAA}$ —

$t_2 = \text{TTTT CGCGTTT}$ —

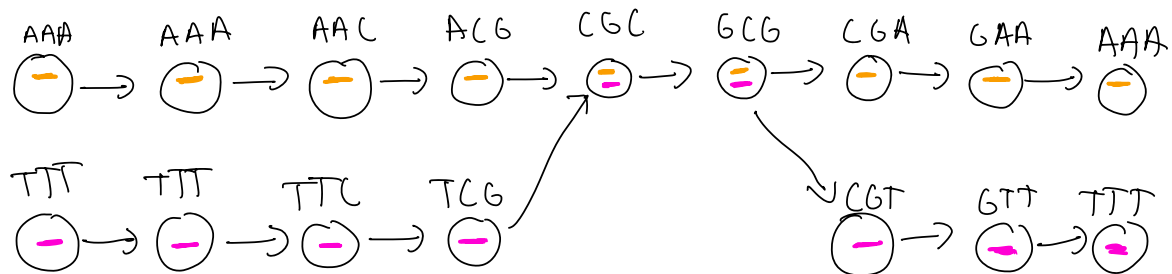
t_1 3-mers

AAA
AAA
AAC
ACG
CGC
GCG
CGA
GAA
AAA

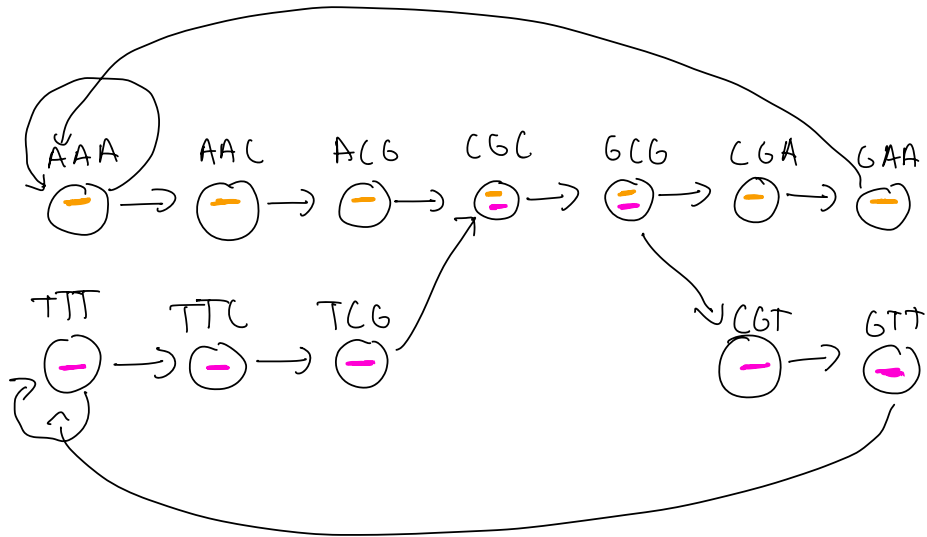
t_2 3-mers

TTT
TTT
TTC
TCG
CGC
GCG
CGT
GTT
TTT

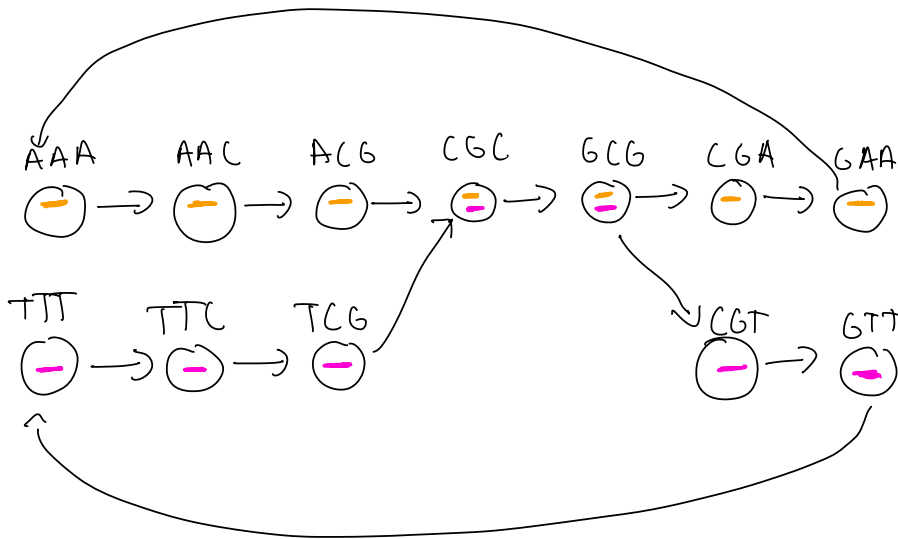
Given the k-mers, I can draw the "nucleotide" graph w/o collapsing.



Best practice is to "compress" the graph. Let's just use multi-edges



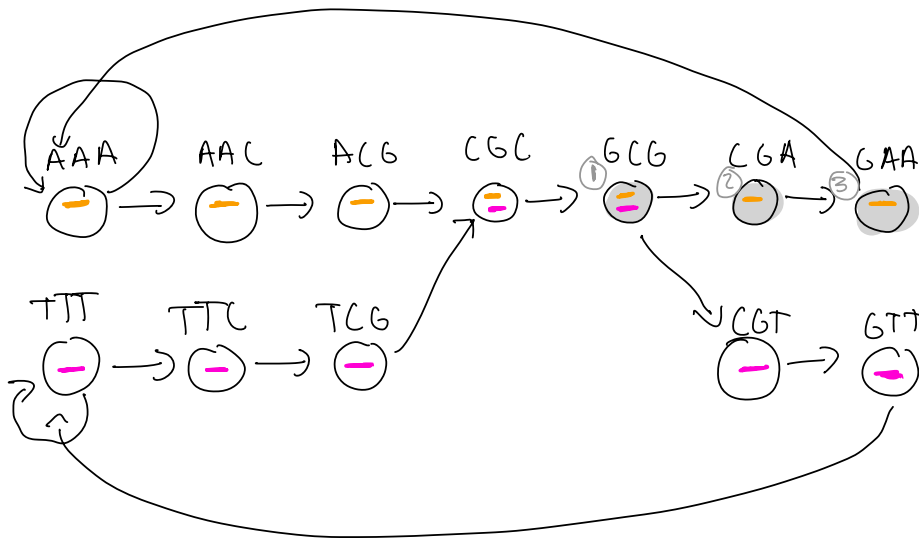
Alternatively:



Problem (2. b)

We will use the multi-edge version from (a) to be a bit more explicit

There are a few possibilities. 1.) There is an error. 2.) your prot. screwed up & meant to give you the sequence GCGAA. Let's work under ass. (2). This case is straightforward as there is one path. Let's number & color the nodes we visit.



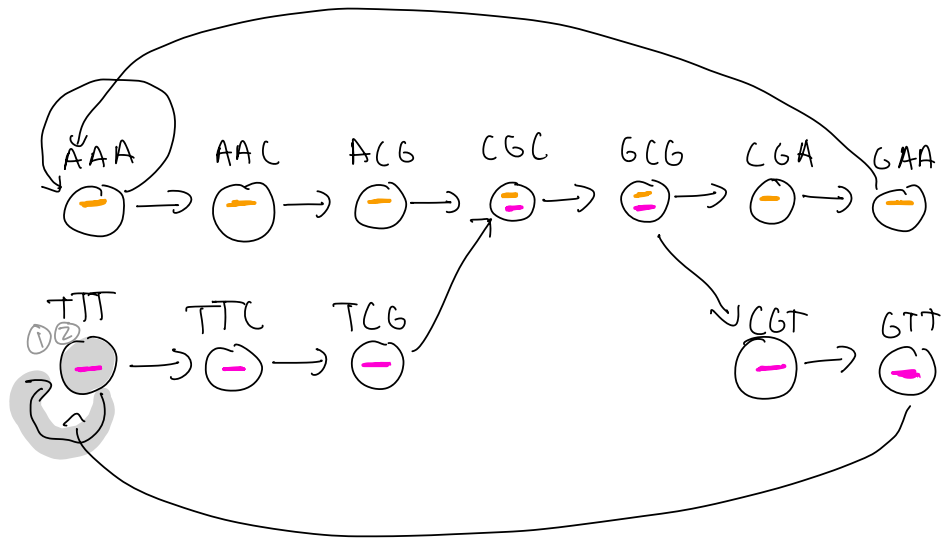
Pseudocode:

$$\text{①} \wedge \text{②} \wedge \text{③} = \text{---}, \text{ so clearly fr.}$$

The problem does not state forward/reverse, so we can pretend they are in the forward strand only as long as we are explicit about this assumption.

Puzzle (2.c)

Because we know all reads come from the forward strand, there is only one possibility.



Pseudalignment:

$$\textcircled{1} \text{ --- } \textcircled{2} = \text{ --- }, \text{ so, } t_2.$$

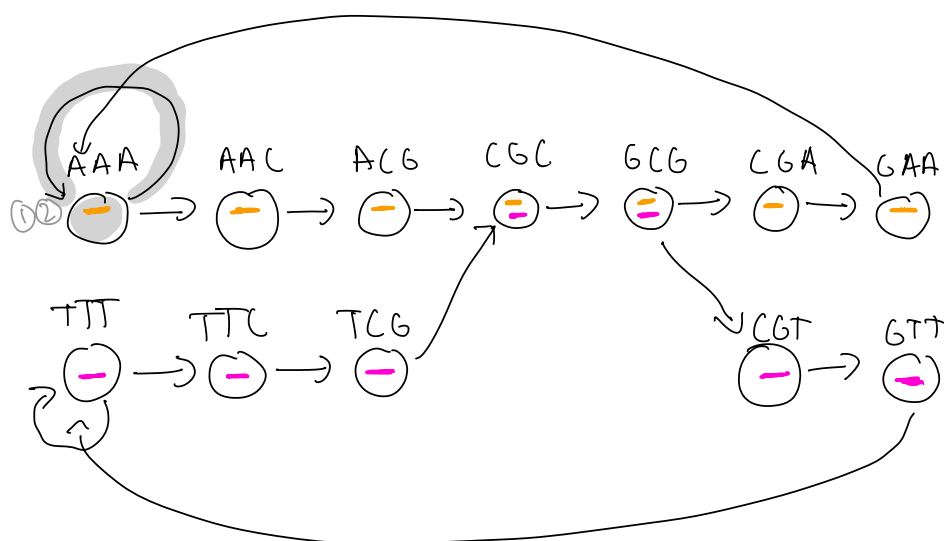
Problem (2d)

There are a few ways to deal with this. See disc. week 6 for graph alternatives. The easiest (& slowest) is to take the reverse complement of the read:

Read: TTTT

Rev. comp: AAAA

The pseudalignment for the direct observation is in (2c). Below is the PA for the rev. complement.



Pseudalignment:

$$\textcircled{1} \text{ --- } \wedge \text{ --- } \textcircled{2} = \text{ --- }$$

Now, we need to "combine" the information somehow. One way is the union:

$$\text{ --- } \cup \text{ --- } = \text{ --- } \text{ Final answer.}$$

We see that this read is ambiguous given the protocol.

Problem (2c)

Again, a few different ways to do this problem.

The key here is to realize that the sequences share many partners because of the reverse complementarity.

In fact, no 4-mer exists that will allow us to uniquely resolve the read:

t_1 : AAAA C G C G AAA

t_1 r.c.: TTT C G C G TTT

t_2 : TTTT C G C G TTT

t_2 r.c.: AAA C G C G AAA

By taking the reverse complement of both transcripts,

we see that every 4-mer of t_1 is in the reverse

complement of t_2 . So, we have shown that there is no

position of length 4 that enables us to uniquely identify

the transcript.