

NOTE TO COMP 462/561 – Fall 2016 students: These are questions taken from the 2014 final exam. It does not include all questions on that exam, because some of those questions covered material we did not cover this year. Consequently, you should expect that your exam will be a bit longer than what you have here.

Question 1. (16 points, 4 points each)

Indicate whether the following statements are true or false. **Give a two line explanation for each.** *Credits will be given only if the justification is clear and correct.*

- a) True or False? *Justify*. Indels that occur in the coding region of a gene will always cause a frame shift.

- b) True or False? *Justify*. An RNA-seq experiment can be used to measure the abundance of each of the 20,000 human proteins in a given sample.

- c) Give *two* reasons why the secondary structure inferred for a given RNA sequence using the Nussinov algorithm may not always reflect the true structure this sequence adopts in cells.

- d) (4 points) Gene expression can be assessed using a RNA-seq experiment. In that case, the expression level of gene *g* is obtained using the FPKM measure (Fragment Per KiloBase Per Million reads).

Explain the two types of normalization this entails, and why they are necessary.

“per KiloBase”:

“per Million Reads”:

Question 2. (16 points)

Suppose that the Philae mission to comet 67P was equipped with a DNA sequencer. It finds DNA on the comet and researchers determine that on that comet, genes have a highly simplified structure:

- Proteins are made of only 4 types of amino acids. Each amino acid is encoded by a single nucleotide.
- Genes always start with ‘ A A ’
- Genes always end with ‘ T T ’
- The body of genes is made of any number of ‘A’, ‘C’, or ‘G’, ‘T’, but never contain two consecutive T’s, as those would be interpreted as a STOP signal. In genes, C’s and G’s are each twice as frequent as A’s and T’s.
- Genes contain no introns
- Genes are on average 100 bp long.

- a) (9 points) Draw an HMM that could be used to make gene predictions in this type of DNA. Include all the transition and emission probabilities. If you think that there’s some information you are missing in order to choose some of these probabilities, mark them as “?”.

- b) (4 points) Indicate what information you would need to be able to choose the value of probabilities you’ve marked as “?”.

- c) (3 points) If no one is able to give you the information you specified in (b), name the algorithm that you could use to estimate them?

Question 3. (15 points)

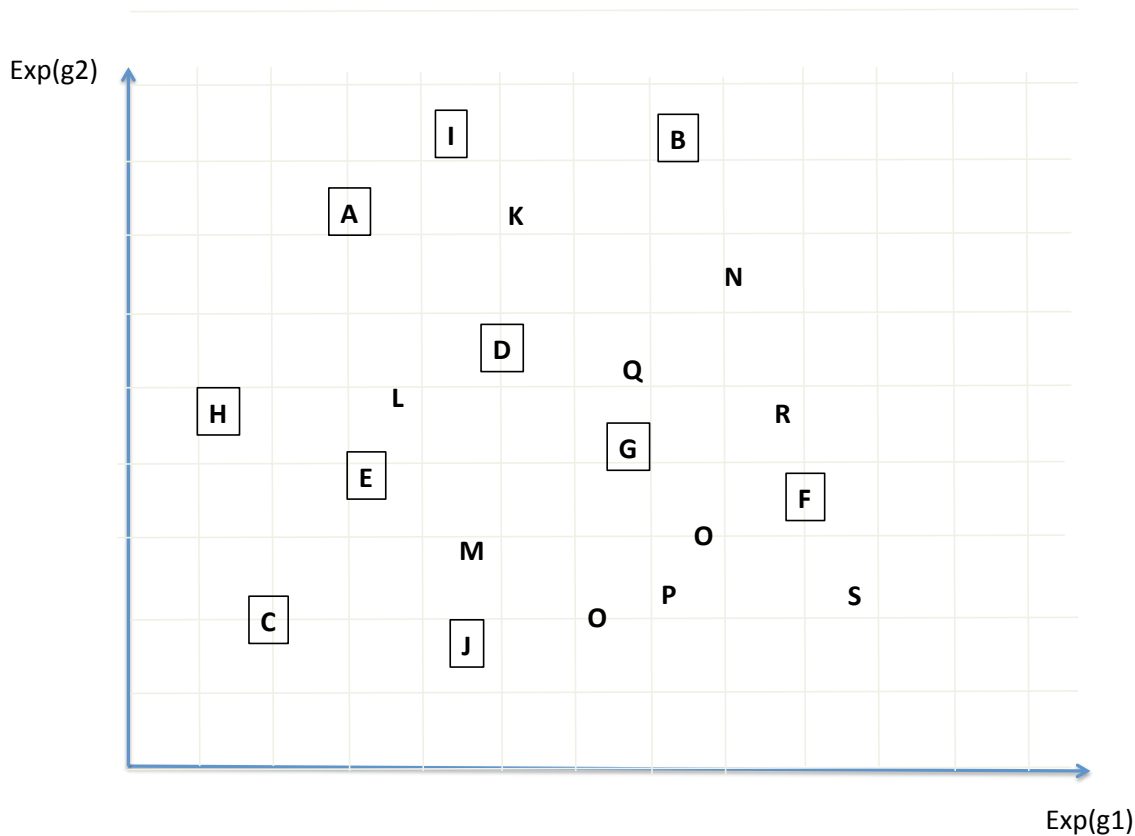
When using Sanger sequencing to read a short piece of DNA, the probability of sequencing error, i.e. calling nucleotide x whereas the correct nucleotide is y , increases as a function of the position within the read. Suppose that a read has length 1000 and that $\Pr[\text{error at position } p] = p / 1000$. Assume that insertions (i.e. the inclusion in a read of a nucleotide that was not present in the DNA sequence) never happen. Deletions (the omission of a nucleotide) are possible but very rare, so we will assume that a read contains at most one deletion of one nucleotide.

a) (10 points) Give a modified version of the Needleman-Wunsch algorithm that could be used to align two reads, under the assumption given below. Pay particular attention to the substitution and indel scoring schemes.

b) (5 points) Suppose that you are interested in reading as accurately as possible the sequence of a portion of a genome and that you have multiple reads originate from that region. Describe informally in 4-5 lines how could you combine the information contained in the reads in order to obtain the most accurate prediction about the exact sequence of that genomic region?

Question 5. (16 points)

Consider the following result of a microarray experiment, where in each case we measured the expression of only two genes, g_1 and g_2 . Assume that samples A, B, ..., J (shown with boxes in the figure below) come from patients with a specific disease, while samples K, L, ..., S (shown without boxes) come from healthy patients.



- (5 points) Draw the linear classifier that obtains the smallest number of classification errors (False-positives + False-negatives) on this data set.
- (3 points) What is the sensitivity of your classifier on this training data?
- (3 points) What is the specificity of your classifier on this training data?
- (5 points) If one wants a sensitivity of at least 90%, where should the boundary of the linear classifier be moved in order to maintain a specificity that remains as high as possible? Draw your answer on the figure above, using a dashed line.

Question 7. (15 Points)

You are given two RNA sequences, $S = s_1 \dots s_m$ and $T = t_1 \dots t_n$, that you want to align. Sequence S has a known nested secondary structure (without pseudoknots), given to you in the form of a list of pairs of positions in S that form base pairs:

$\text{Struct} = \{ (l_1, r_1), (l_2, r_2), \dots, (l_k, r_k) \}$, where $1 \leq l_i < r_i \leq m$ for all $i \in \{1 \dots k\}$.

The secondary structure of sequence T is not known but is believed to be related to that of S .

Give an alignment algorithm to align sequences S and T using a given substitution matrix M and linear gap penalty c , subject to the following constraint:

If nucleotides s_i and s_j form a base pair in S (i.e. $(i, j) \in \text{Struct}$), and t_a is aligned to s_i and t_b is aligned to s_j , then t_a and t_b must be complementary nucleotides (i.e. A-U, C-G, G-C, or U-A pairs).

For example, if S and T are the sequence shown below and S has the structure that is shown,

$S = A \ C \ G \ U \ G \ A \ A \ A \ C \ C \ G \ U$

$T = A \ G \ A \ A \ C \ A \ U \ C \ C$

Then the following alignment is admissible:

$S = A \ C \ G \ U \ G \ A \ A \ A \ C \ C \ G \ U$

$T = A \ G \ - \ - \ A \ A \ C \ A \ U \ C \ C \ -$

But this one is not:

$S = A \ C \ G \ U \ G \ A \ A \ A \ C \ C \ G \ U$

$T = - \ - \ A \ G \ A \ A \ C \ A \ U \ C \ C \ -$

Page left blank intentionally. Use it if you need extra space.

Page left blank intentionally. Use it if you need extra space.