

# Examination Bioinformatics

Linköpings Universitet, IDA, Statistik

---

Course:	732A51 Bioinformatics
Date:	2021/02/08, 8–12
Teacher:	Krzysztof Bartoszek
Provided aids:	The help material is included in the zip file <b>exam_help_material_732A51.zip</b> .
Grades:	A= [18 – 20] points B= [16 – 18) points C= [14 – 16) points D= [12 – 14) points E= [10 – 12) points F= [0 – 10) points
Instructions:	<p>Provide a detailed report that includes plots, conclusions and interpretations. Give motivated answers to the questions. If an answer is not motivated, the points are reduced. Provide all necessary codes in an appendix.</p> <p>If you are asked to do plots, then make sure that they are informative, have correctly labelled axes, informative axes limits and are correctly described.</p> <p>Points may be deducted for poorly done graphs.</p> <p>Name your digital part solution files as: <b>[your id]_[own file description].[format]</b></p> <p>If you have problems with creating a pdf you may submit your solutions in text files with unambiguous references to graphics and code that are saved in separate files</p> <p>There are <b>THREE</b> assignments (with sub-questions) to solve.</p> <p>Include all code that was used to obtain your answers in your solution files.</p> <p>Make sure it is clear which code section corresponds to which question.</p> <p>Your code should be complete and readable, possible to run by copying directly into a script. Comment directly in the code whenever something needs to be explained or discussed.</p> <p>If you also need to provide some hand-written derivations please number each page according to the pattern: Question number . page in question number i.e. Q1.1, Q1.2, Q1.3, ..., Q2.1, Q2.2, ..., Q3.1, ... .</p> <p>Scan/take photos of such derivations preferably into a single pdf file but if this is not possible multiple pdf or .bmp/.jpg/.png files are fine.</p> <p>Please do not use other formats for scanned/photographed solutions.</p> <p>Please submit all your solutions via LISAM or e-mail. If emailing, please email them to <b>BOTH</b> krzysztof.bartoszek@liu.se and KB_LiU_exam@protonmail.ch .</p> <p>During the exam you may ask the examiner questions by emailing them to KB_LiU_exam@protonmail.ch <b>ONLY</b>. Other exam procedures in LISAM.</p>

---

## Problem 1 (7p)

A) (5p) The overall base composition of a new (in 2006) bacterial strain of the *Chlorobium* genus is G=C= 46.7% and A=T= 53.3%. In a random (all positions independent and identically distributed) sequence of 2,572,079 bp nucleotides with these proportions, what is the expected number of occurrences of the sequences CTAG, TCAG and GCTA?

B) (2p) The number 2,572,079 bp is A) is not taken randomly. Find how it is related to the problem. You will have to do an internet search (you are permitted for this question to search anywhere online during the exam). A starting point for the online search can be the article describing the new strain, [VoglGlaeserPfannesWannerOvermann\\_Chlorobium\\_2006.pdf](#). Do **not** read it in detail during the exam, think what can be the key piece of information for finding the origin of this number. Write (very briefly) how this number is related to the bacterial strain in question.

## Problem 2 (8p)

For each matrix A) PAM250 (file [PAM250.png](#)) and B) BLOSUM62 (file [BLOSUM62.png](#)) answer the following

1. What are/is the most probable substitution? If there are too many to write out, then provide the value in the table and three examples.
2. What are/is the least probable substitution? If there are too many to write out, then provide the value in the table and three examples.
3. Which substitution is more probable  $W \rightarrow F$  or  $H \rightarrow R$ ?
4. Which substitution is more probable  $M \rightarrow Q$  or  $K \rightarrow T$ ?

**TIP:** The numbers in the matrices are some transformations of the original observed data. Exactly how this transformation was done and what do the actual numbers mean is not that important for this question.

## Problem 3 (5p)

A) (3p) You are given, in Newick format, the following two rooted trees with tip labels

$$T_1 = (A, (((E, D), C), ((G, F), (B, H))))$$

and

$$T_2 = (((((H, B), (F, G)), (C, (D, E))), A).$$

Are they identical in topology?

B) (2p) Draw all possible rooted, without tip labels, trees with 3, 4 and 5 tips.