

CENG 465 – Introduction to Bioinformatics

Spring 2019-2020

Assignment #4: Analysis of Microarray Data

Due Date: June 7, 2020, Sunday, 11:59 PM

Data Analysis Assignment about Microarrays

In this assignment, you are going to analyze a microarray dataset of 30 samples. The experiment is performed on the human genome and contains gene expression levels as signal intensities for 22215 human transcripts, i.e., mRNAs that encode proteins. The experiments are performed on two different types of samples: disease and healthy tissues. In other words m samples are taken from diseased tissues and n samples are taken from healthy tissues. $m+n = 30$, and m may or may not be equal to n . You have two main goals in this assignment:

Goal 1: Determine which sample is taken from which tissue. What is m ? What is n ? Provide a matching for each sample to one of the two tissue types. The labels of the tissues can be chosen arbitrarily, i.e., it is OK if you swap the correct labels. Here, the goal is to correctly cluster the samples.

Goal 2: Find 10 genes each that are expressed specifically in each of the tissues. In other words find 10 genes that are highly expressed in the **diseased** tissue and not expressed in the **healthy** tissue. Similarly find 10 genes that are highly expressed in the **healthy** tissue and not expressed in the **diseased** tissue. You should list a total of 20 genes.

You are free to use any of the techniques we have learned in class, or techniques you already know, or techniques you have just invented. You may write a program or use existing tools. You should clearly describe the methods you have used to accomplish the goals given above. One example tool that you can use is TM4 MEV (<http://mev.tm4.org/#/welcome>). You may use your own judgment for any issue that is not specified clearly in this text.

The Dataset:

- The complete dataset can be downloaded as a single tab separated text file from the following address:
 - <http://user.ceng.metu.edu.tr/~tcan/download/hw4dataset.txt>

Deliverables:

- A short report which contains a step by step description of the tasks that you have performed, the grouping of the samples into the two tissues, and the list of 10+10 genes specifically expressed in each tissue.

Submission: Submit the deliverable as a .txt, PDF, or Word document via ODTU-Class.

Late Submission Policy: Your final assignment grade will be penalized 20 points per late day.