## GS01 0163 Assignment 3 Fall 2007 Due Date: 4 October 2007

**Metaquestion 0**. We discussed in class how Sweave could be used to document the use of R code by combining code and documentation in a single report. Examples are available on the Lecture Notes page for lectures 6 and 7. For this assignment, which involves coding in R, we would like you to prepare your reports using Sweave. Submissions would ideally involve two files: (a) the source file, which we should be able to run on our own machines, and (b) the final output pdf file, which should include the results of running the R code. Sweave itself is already available as part of R. In order to produce the pdf output, however, you will need MiKTeX as well (http://miktex.org/). Once MiKTeX is installed, you can (a) prepare the template Sweave document (e.g., "fileA.Rnw"), (b) invoke Sweave('fileA.Rnw'), which will produce fileA.tex, and (c) from the windows command prompt in the appropriate directory, invoke "pdflatex fileA.tex", which should produce fileA.pdf.

- 1. This assignment uses a subset of the Singh prostate cancer data that was acquired for Homework 1. Download the "subcel.txt" and "subsamples.txt" files for homework 3 from the course web site. These files list the 20 CEL files to be used in this assignment. (You will have to edit "subcel.txt" so the path points to the correct location on your computer.)
  - (a) Use the supplied "subsamples.txt" file to create a phenoData object. Use short sample names to identify the arrays.
  - (b) Load in all 20 CEL files to create an AffyBatch object. Create a boxplot and a histogram of the raw data. Do the arrays need to be normalized?
  - (c) Look at images of several of the arrays. Do you see anything unusual?
- 2. Use BioConductor to decide if the RNA used in any of these experiments was exceptionally degraded.

Bonus: What happens if you plot the degradation slopes as a function of the factors in the pData object in the AffyBatch?

- 3. In Homework 1, we saw that clustering the samples after standard processing with dChip, there were two dominant clusters that did not match the biological division into normal and cancer. The 20 samples used in this homework are evenly balanced between normal/cancer and between the two dChip clusters. Using the qc function in the simpleaffy package, explore the various quality control metrics. Do any of the quality control metrics suggest an explanation for the differences between the two clusters?
- 4. Perform background correction on all the arrays using two different methods: "mas" and "rma". Prepare histograms and boxplots of the AffyBatch objects for both background correction methods. What can you say about the differences? Which background correction method do you think is better? Why?

- 5. This problem is a continuation of the previous problem. This problem, like several of the later problems, uses the simpleCluster function that is available from the course web site. This function produces cluster dendrograms based on the 25% of genes with the highest average expression.
  - (a) After processing the arrays using "mas" background correction, use expresso to perform "quantiles" normalization, select the "pmonly" features, and quantify using the "medianpolish" summarization method. Use simpleCluster to produce a dendrogram.
  - (b) Repeat part (a) using "rma" background correction.
  - (c) Does the background correction method have any effect on clustering based on the highest expressing genes?