



Northeastern University

College of Science

Module 11 – Homework

Problem 1: (40 points)

Clustering analysis on the "CCND3 Cyclin D3" gene expression values of the Golub et al. (1999) data.

- (a) Conduct hierarchical clustering using single linkage and Ward linkage. Plot the cluster dendrogram for both fit. Get two clusters from each of the methods. Use function `table()` to compare the clusters with the two patient groups ALL/AML. Which linkage function seems to work better here?
- (b) Use *k*-means cluster analysis to get two clusters. Use `table()` to compare the two clusters with the two patient groups ALL/AML.
- (c) Which clustering approach (hierarchical versus *k*-means) produce the best matches to the two diagnose groups ALL/AML?
- (d) Find the two cluster means from the *k*-means cluster analysis. Perform a bootstrap on the cluster means. Do the confidence intervals for the cluster means overlap? Which of these two cluster means is estimated more accurately?
- (e) Produce a plot of *K* versus SSE, for *K*=1, ..., 30. How many clusters does this plot suggest?



Northeastern University

College of Science

Problem 2 (30 points):

Cluster analysis on part of Golub data.

- (a) Select the oncogenes and antigens from the Golub data. (Hint: Use `grep()`).
- (b) On the selected data, do clustering analysis for the genes (not for the patients). Using K-means and K-medoids with $K=2$ to cluster the genes. Use `table()` to compare the resulting two clusters with the two gene groups oncogenes and antigens for each of the two clustering analysis.
- (c) Use appropriate tests (from previous modules) to test the marginal independence in the two by two tables in (b). Which clustering method provides clusters related to the two gene groups?
- (d) Plot the cluster dendrograms for this part of golub data with single linkage and complete linkage, using Euclidean distance.



Northeastern University

College of Science

Problem 3 (30 points):

Clustering analysis on NCI60 cancer cell line microarray data (Ross et al. 2000)

We use the data set in package ISLR from r-project (Not Bioconductor). You can use the following commands to load the data set.

```
install.packages('ISLR')  
library(ISLR)  
ncidata<-NCI60$data  
ncilabs<-NCI60$labs
```

The `ncidata` (64 by 6830 matrix) contains 6830 gene expression measurements on 64 cancer cell lines. The cancer cell lines labels are contained in `ncilabs`. We do clustering analysis on the 64 cell lines (the rows).

- (a) Using k-means clustering, produce a plot of K versus SSE, for $K=1, \dots, 30$. How many clusters appears to be there?
- (b) Do K-medoids clustering ($K=7$) with 1-correlation as the dissimilarity measure on the data. Compare the clusters with the cell lines. Which types of cancer are well identified in a cluster? Which types of cancer are not grouped into a cluster? According to the clustering results, which types of cancer are most similar to ovarian cancer?