

**FES758b / S&DS363 / S&DS 563**  
**Multivariate Statistics**  
**Homework #3: Cluster Analysis**  
Due: Friday, 2/22/2019 11:59pm on CANVAS

*Answers should be complete and concise. You may use any statistics program for calculations that you wish (multiple programs is fine too).*

**IF YOU WORK IN A GROUP, YOU MAY TURN IN ONE ASSIGNMENT FOR YOUR GROUP. BE SURE ALL NAMES ARE ON THE ASSIGNMENT.**

**Examples of Cluster Analysis in SAS and R are up on the  
CANVAS home page**

## SAMPLE DATA SET

**The example below is JUST FOR YOUR PRACTICE.  
NOTHING TO TURN IN HERE!**

The file senate104.xls (or senate104.csv) contains the voting records for US senators during the 104<sup>th</sup> congressional session. This file actually only contains the first 198 votes. For each vote, senators may respond yes (1) or no(0). They may also abstain or not be present. For this dataset, I have replaced abstain/not present votes with 0.5 (hierarchical clustering in many software packages will not work with missing values). Your task is to cluster senators.

Note : there are more than 100 senators because some senators resigned during the session and were replaced by others. Also, note that one Colorado senator switched from being a democrat to being a republican and is thus listed twice.

1). Get a measure of the standard deviation of each vote. What do you observe? (no need to turn in all the output). Votes with low standard deviation will not have much affect on the clustering process (i.e. no differentiation between senators based on these votes).

*Note : in SPSS, use Analyze → Descriptive Statistics → Descriptives. In SAS, use something like*

```
proc means data=in.senate104;  
  var v1-v198;  
run;
```

*In R, assuming the votes are in an object called senate, use*

```
#get the data from a CSV file
senate=read.csv("http://reuningscherer.net/stat660/data/senate104.csv",header=T)
```

```
#get standard deviation for each vote
round(sqrt(apply(senate[, -1], 2, var)), 2)
```

*Standard deviations were largely in the .3, .4 range; however, some are lower in the .1 and less range – that is, some votes had more agreement.*

2). Use hierarchical clustering on the data. Try two metrics and two agglomeration procedures, one of which should be complete linkage (furthest neighbor). You might try some measure appropriate to binary data (which this is approximately).

*Note – in SAS, I've given some code below. By default, Proc Cluster only calculates Euclidean distance. However, you can use a macro called 'distance' to calculate other distances. Copy the file distnew.sas to your computer from the Software folder on the classes server in the materials folder.*

*The usage is as follows :*

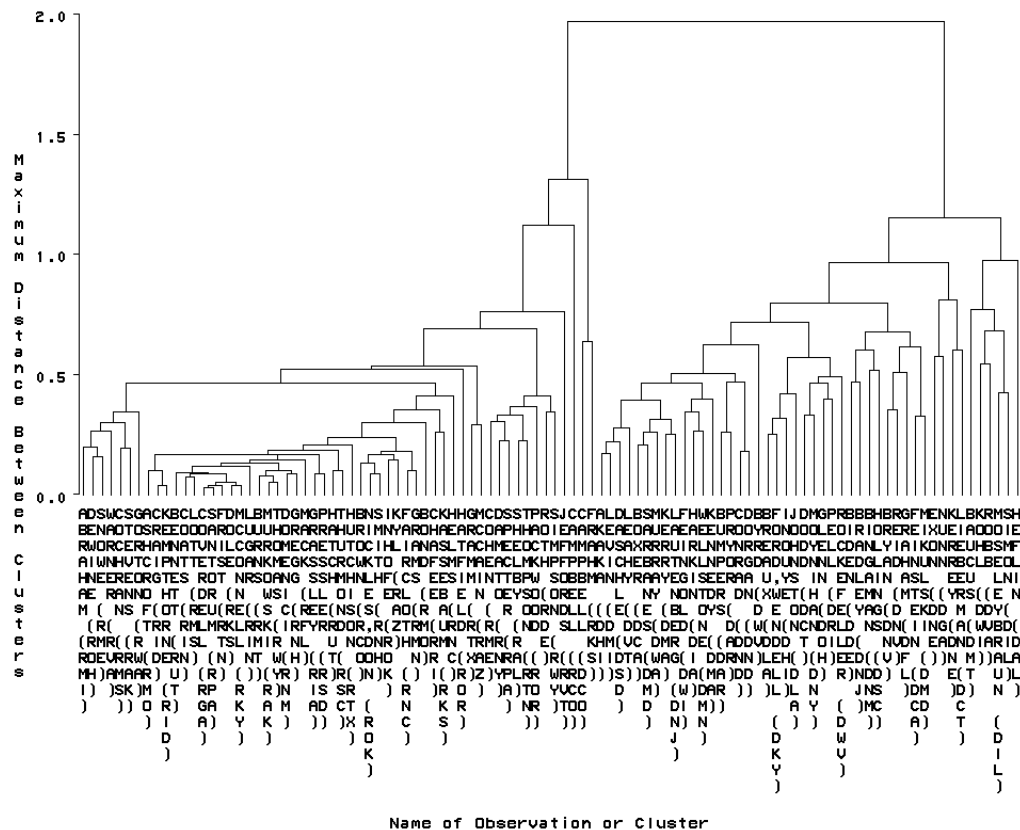
```
*change method= for other metrics. See SAS help file for options.
*Use SNORM option to standardize variables. Need to use other options
*in the VAR statement for categorical data;
```

```
PROC DISTANCE DATA= MYLIB.senate OUT=OUTDIST METHOD=CITYBLOCK
SNORM;
  VAR INTERVAL (V1-V198);
  COPY CEREAL;
RUN;
```

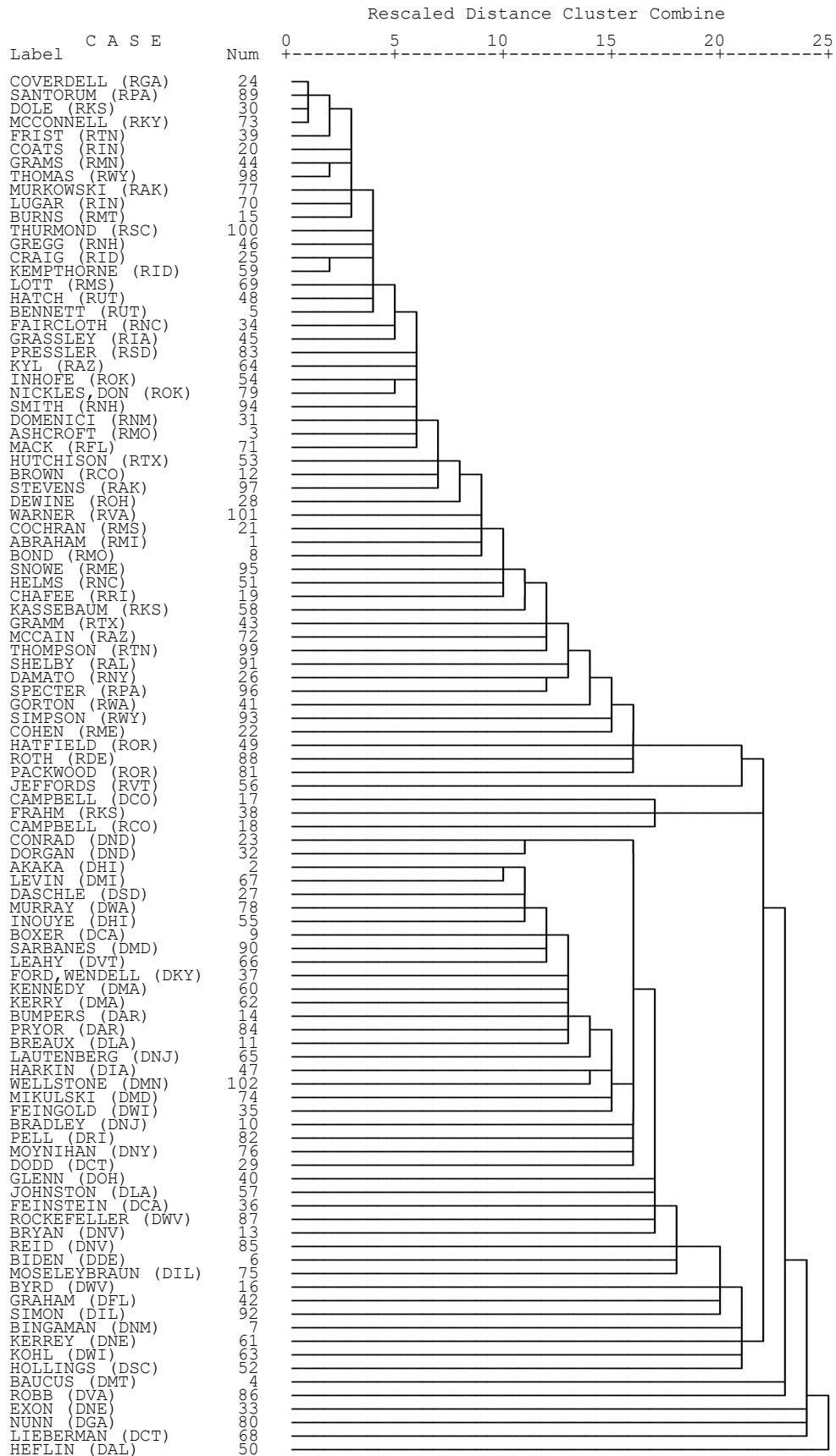
```
*RUN CLUSTERING PROCEDURE AND SAVE RESULTS;
PROC CLUSTER DATA=outdist METHOD=compact RMSSTD RSQ OUTTREE=TREE;
id senator;
RUN;
```

```
*MAKE A DENDOGRAM;
PROC TREE DATA=TREE;
RUN;
```

*Obviously, there are many possibilities here. Here is compact clustering using a manhattan metric (from SAS – see program above).*

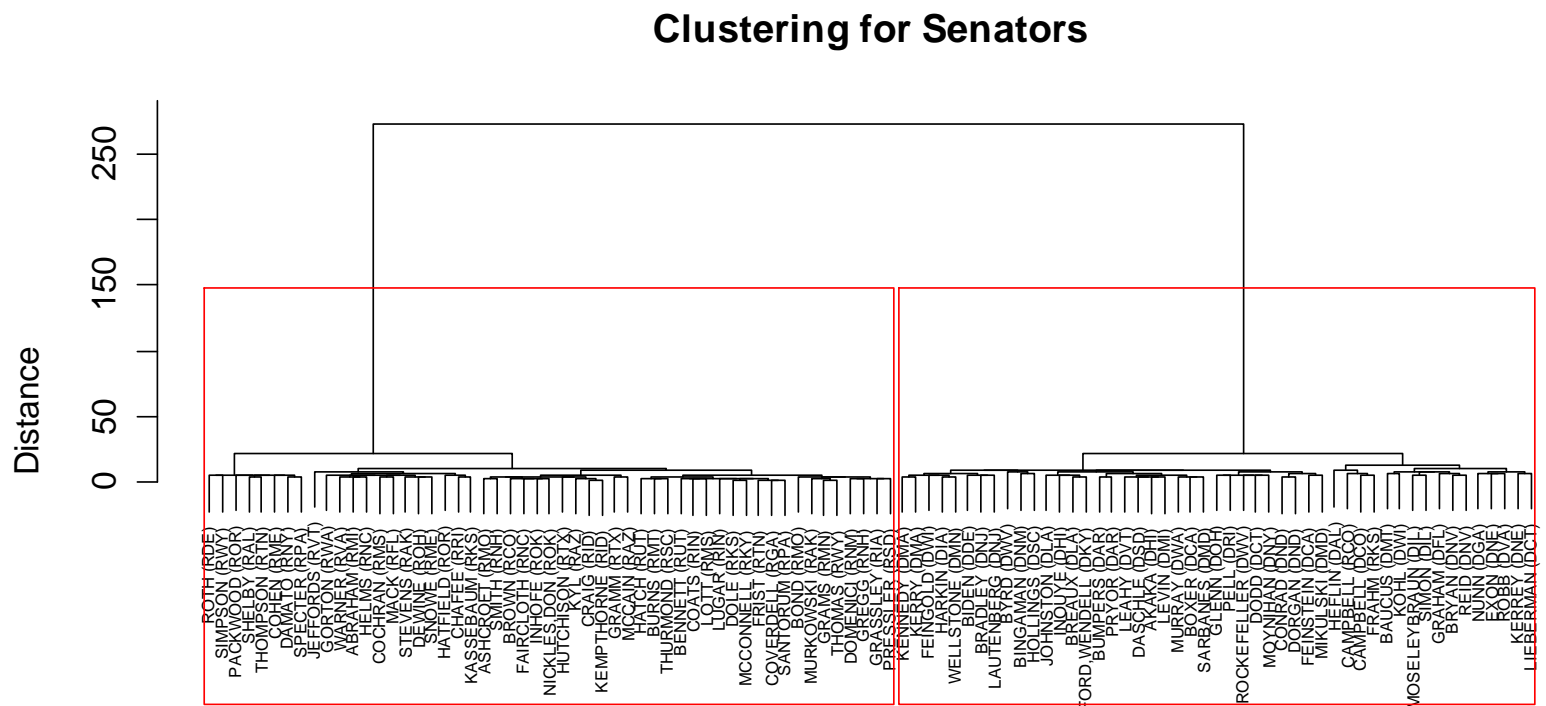


Here is single linkage using using Euclidean distance in SPSS : (Note that this is just a text picture, so you can change font size (below is 8) and change line spacing (below is set to 6 point) to make the picture fit better.



In R, use something like the following :

```
sennorm<-scale(senate[, -1])
#get the distance matrix
dist1<-dist(sennorm [, -1],method="euclidean")
#now do clustering;
clust1<-hclust(dist1,method="ward")
#draw the dendrogram
plot(clust1,labels=senate[, 1],cex=0.5,xlab="",ylab="Distance"
,main="Clustering for Senators")
rect.hclust(clust1,k=2)
```



3). Below is a print of four measures evaluating hierarchical clustering using Euclidean distance and complete linkage (I'm not saying this is the best hierarchical clustering method to use, just the most common). Comment on these four statistics and decide on a suggested number of clusters.

*Here is the SAS Code used to get the plot below. NOTE THAT YOU CANNOT GET THIS PLOT IF YOU START BY USING THE DISTANCE MACRO (I.E. YOU CAN ONLY GET THE PLOT BELOW USING EUCLIDEAN DISTANCE – SEEMS SILLY, BUT I HAVEN'T WORKED OUT A WAY AROUND THIS!!)*

```
PROC CLUSTER DATA=IN.SENATE104 METHOD=compact RMSSTD RSQ OUTTREE=TREE;
id senator;
RUN;
```

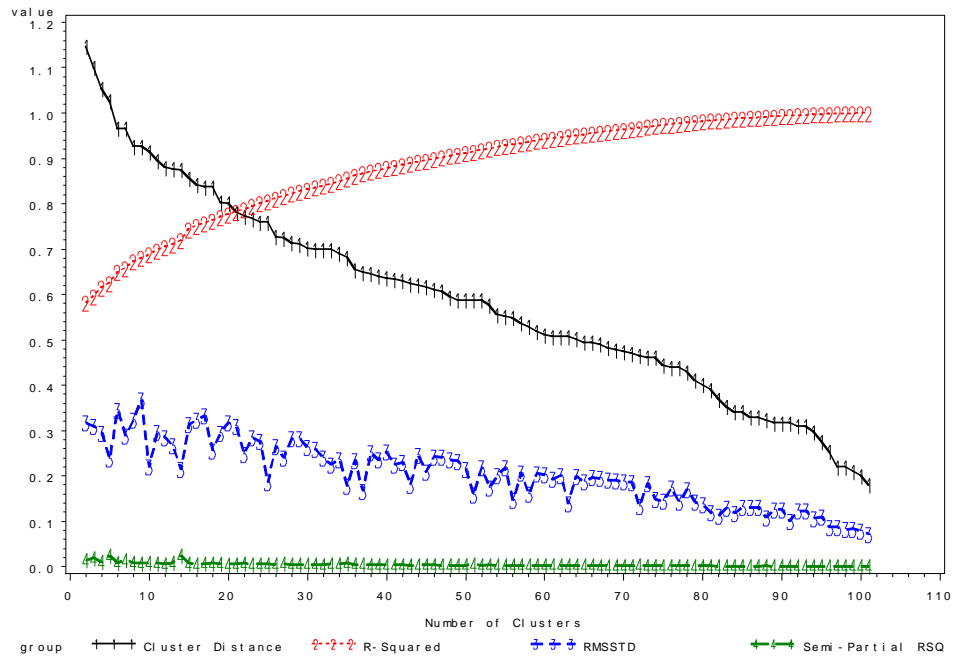
```

*INCLUDE THE FILE WITH THE MACRO FOR EVALUATING CLUSTER NUMBER;
*NEED TO CHANGE FILE LOCATION TO MATCH THAT ON YOUR COMPUTER!!!!;
%INCLUDE 'C:\CLUSTER.SAS';
*RUN THE MACRO - ONLY ARGUMENT IS THE NAME OF THE OUTPUT DATASET FROM
THE CLUSTERING PROCEDURE;
%CLUSTPLOT(TREE);
RUN;

```

*There is no obvious number of clusters suggested by these metrics. Cluster distance perhaps suggest a number of clusters in the 5 to 7 range, but this is not a well-defined break point.*

4). Use k-means clustering on the data. Try somewhere between 2 and 10 groups. Make a scree plot of internal SS versus k to look for an elbow.



## **HOMEWORK ASSIGNMENT**

**PLEASE turn in the following answers for YOUR DATASET! If Cluster Analysis is not appropriate for your data, use one of the two loaner datasets described at the end of the assignment.**

**List your Name(s – if a group) and a one sentence reminder of which dataset your are using.**

1. Think about what metrics are appropriate for your data based on data type. Write a few sentences about this. Also think about whether you should standardize or transform your data (comment as appropriate).
2. Try various forms of hierarchical cluster analysis. Try at least two different metrics and two agglomeration methods. Produce dendrograms and comment on what you observe.
3. If possible, run the SAS macro or the R function to think about how many groups you want to retain (i.e. get the plots of cluster distance, R-squared, etc). If you can't run this, discuss how many groups you think are present.
4. Run k-means clustering on your data. Compare results to what you got in 3.) Include a sum of squares vs. k (number of clusters) plot and comment on how many groups exist.
5. Comment on the number of groups that seem to be present based on what you find above.

# LOANER DATASET 1

(if Cluster Analysis is not appropriate for your data)

The file `stream.xls` contains data on the prevalence of 11 species of microcrustacea in seven streams in Alaska. Five measurements were made at each site. The species are

sp1	<i>Nitocra hibernica</i>
sp2	<i>Atheyella illinoisensis</i>
sp3	<i>Atheyella idahoensis</i> ,
sp4	<i>Bryocamptus hiemalis</i>
sp5	<i>Bryocamptus zschokkei</i> ,
sp6	<i>Acanthocyclops vernalis</i>
sp7	<i>Alona guttata</i> ,
sp8	<i>Graptoleberis</i>
sp9	<i>Chydorus</i>
sp10	<i>macrothricidae</i>
sp11	<i>Maraenobiotus insegmentus</i>

Stream age, the Pfrank index (a measure of stability – higher values are lower stability), Temperature in degrees C, turbidity, conductivity, and alkalinity were also measured for each stream – these values are also included.

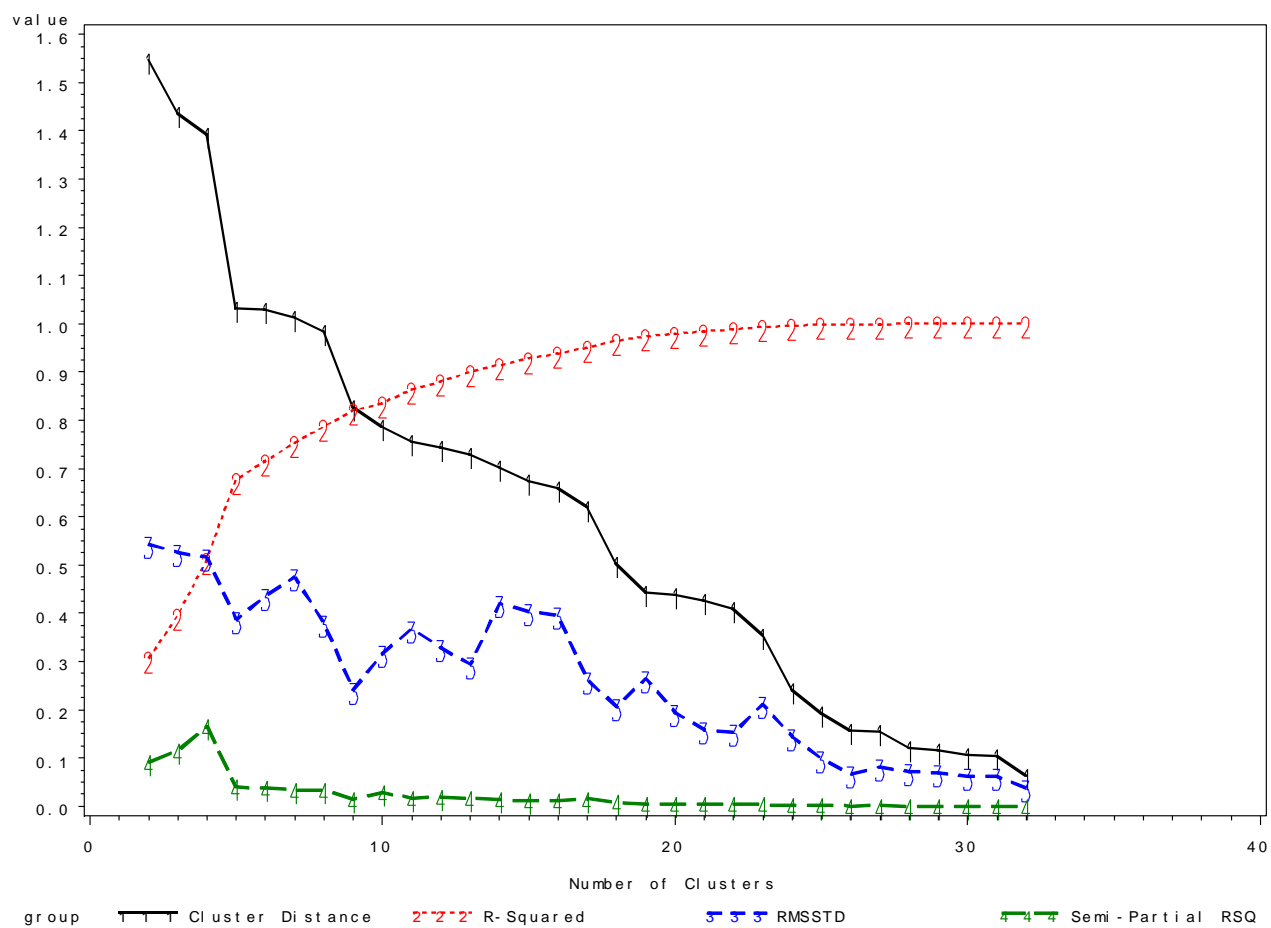
Data is taken from Peter Shaw's Multivariate Statistics for the Environmental Sciences (2003), p. 24

Your goal is to use cluster analysis to identify similar streams.

You'll want to make a transformation of the species data before clustering. I recommend a log transformation, but you'll need to add one before taking logs (because . . . . *can't take log of zero!*)

Below is a SAS printout of statistics used to evaluate the number of clusters based on Euclidean distance and Ward's Method. How many clusters do you think exist?





## **LOANER DATASET 2**

**(if Cluster Analysis is not appropriate for your data)**

The file `University.csv` contains data on from 1995 on 25 Universities. The variables are

- SAT Score
- Percent of Class in Top 10% of high school class
- Acceptance Rate
- Student/Faculty Ratio
- Expenses (dollars)
- Graduation Rate (%)

Use cluster analysis to find groups of Universities (follow instructions for other datasets)